

Наша Машина баз данных (как Oracle Exadata, только про PostgreSQL) и система управления к ней

Константин Аристов,
Скала^р, техлид



HighLoad⁺⁺
2022

Яндекс

Скала^р сегодня:



разработка и производство модульной платформы для высоконагруженных государственных и корпоративных информационных систем

6 лет

серийного
выпуска

150+

комплексов
в промышленной
эксплуатации

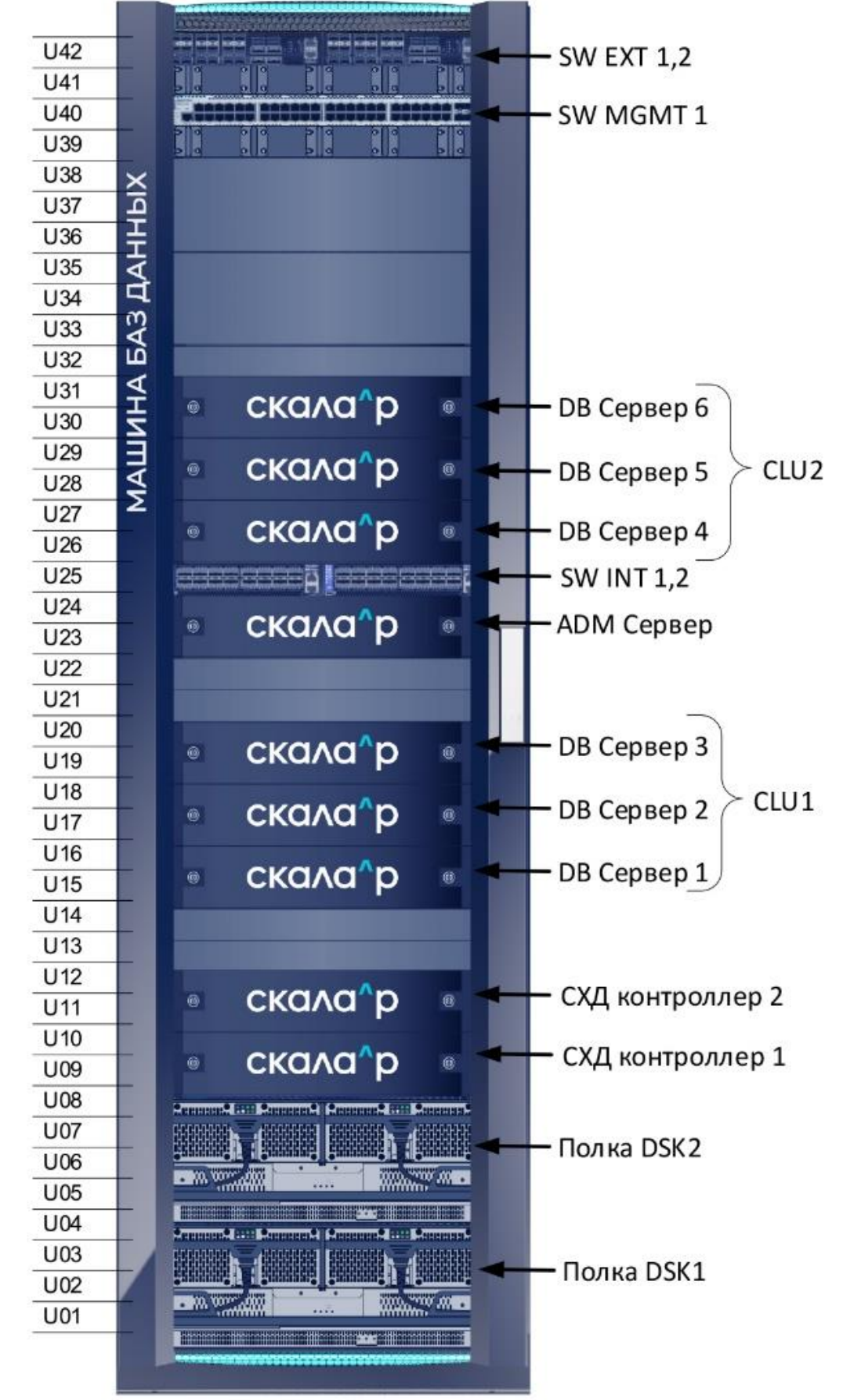
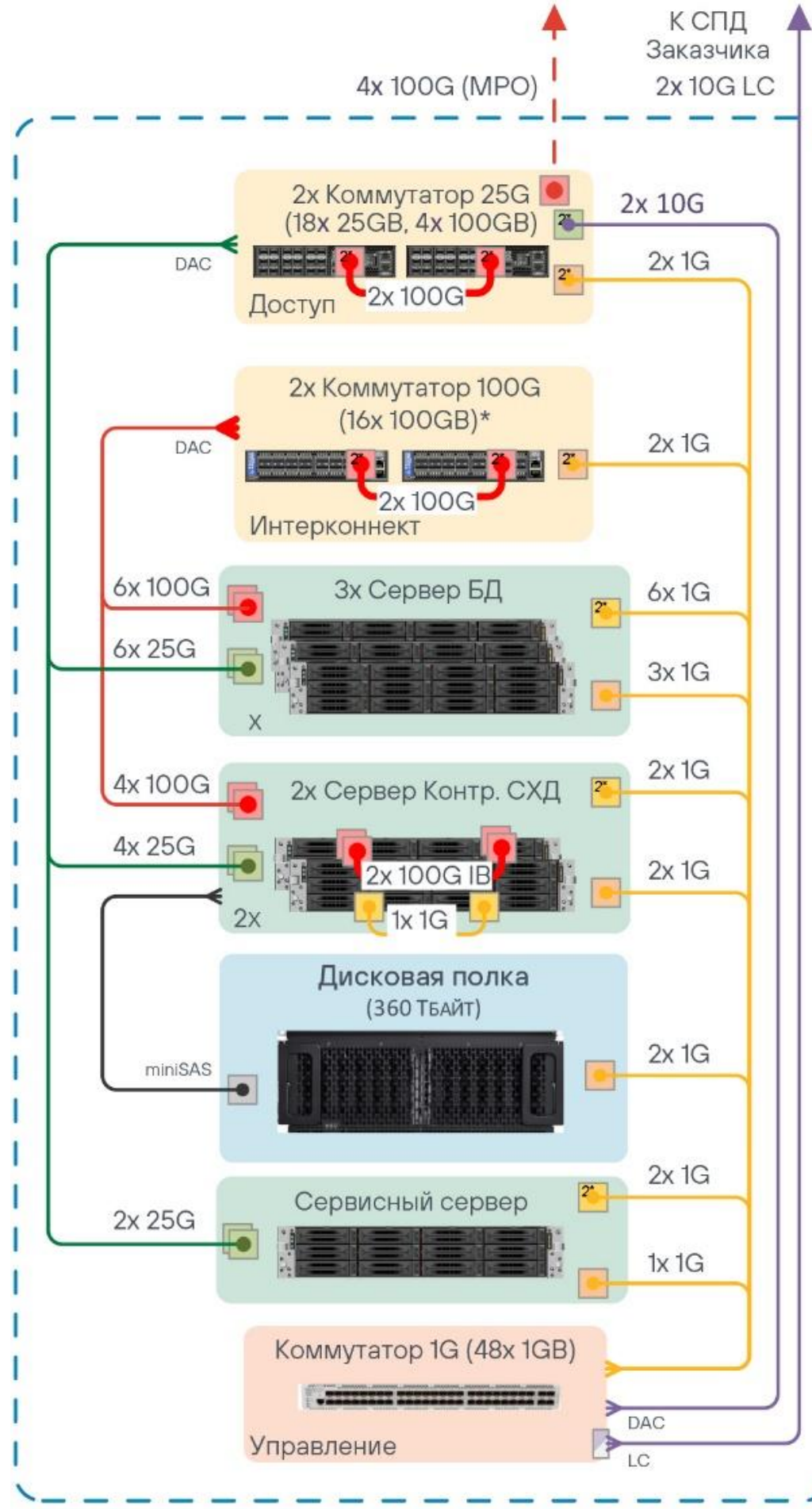
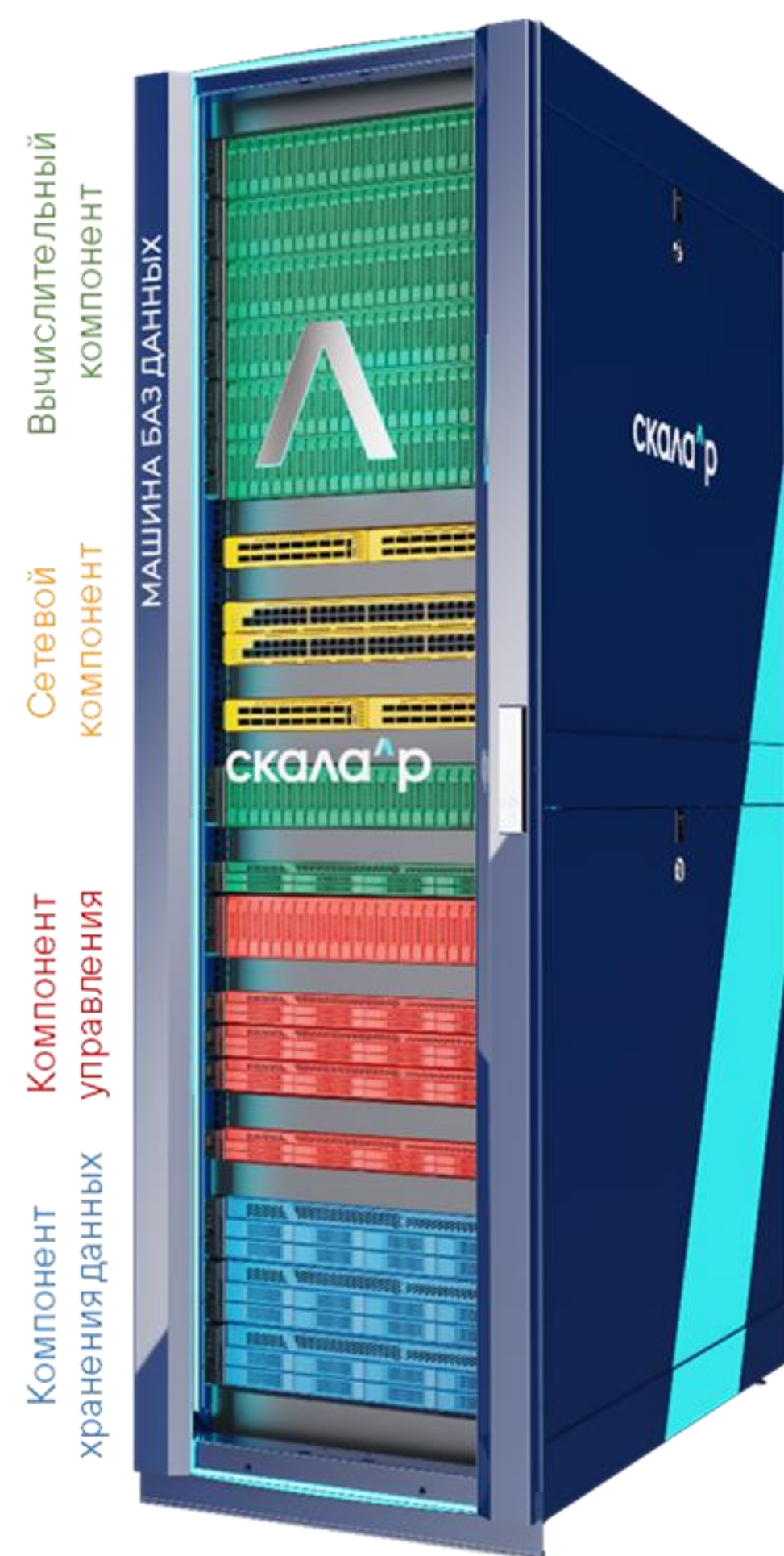
2500+

вычислительных
узлов

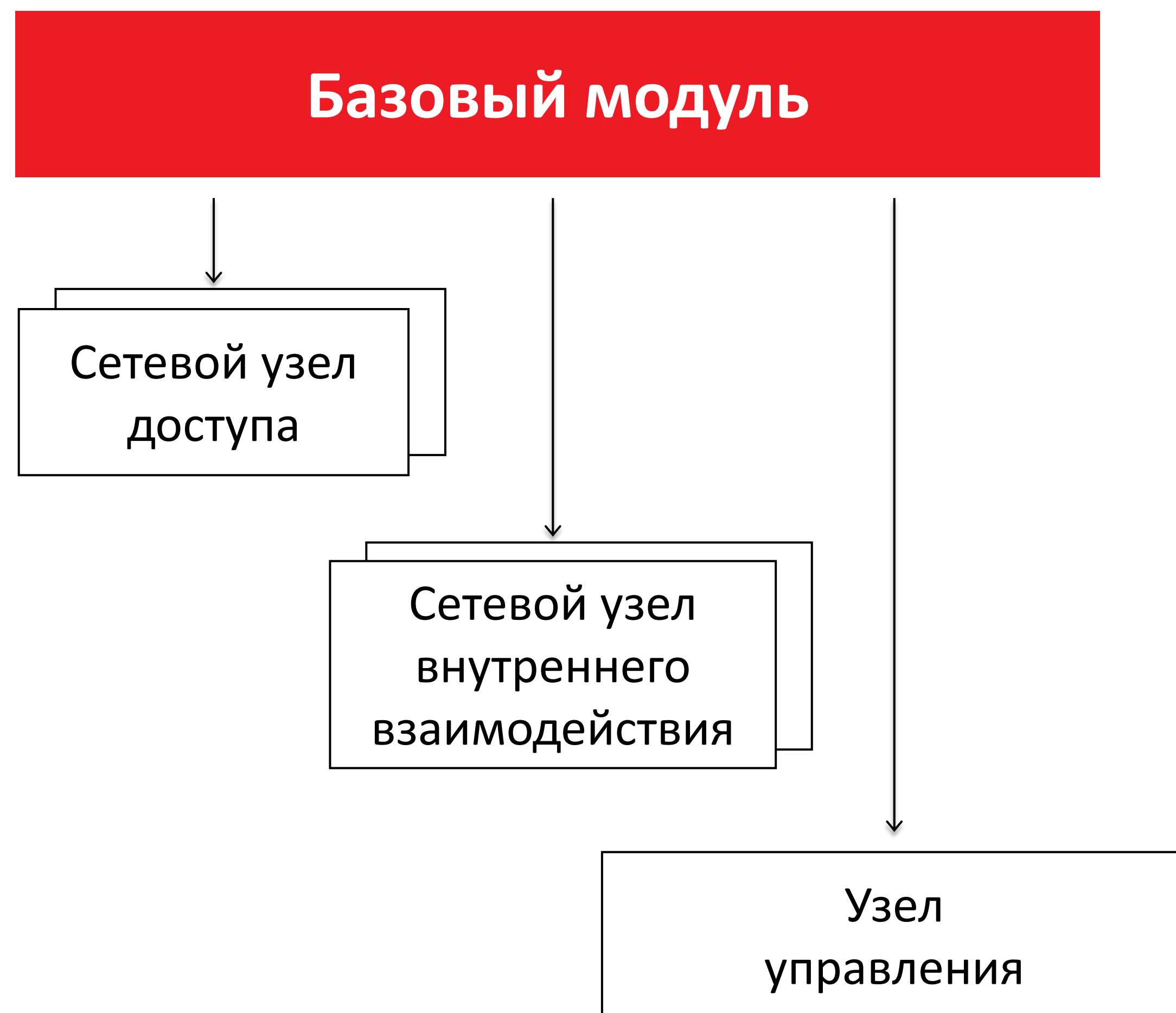
Яндекс



Снаружи и внутри



Снаружи и внутри



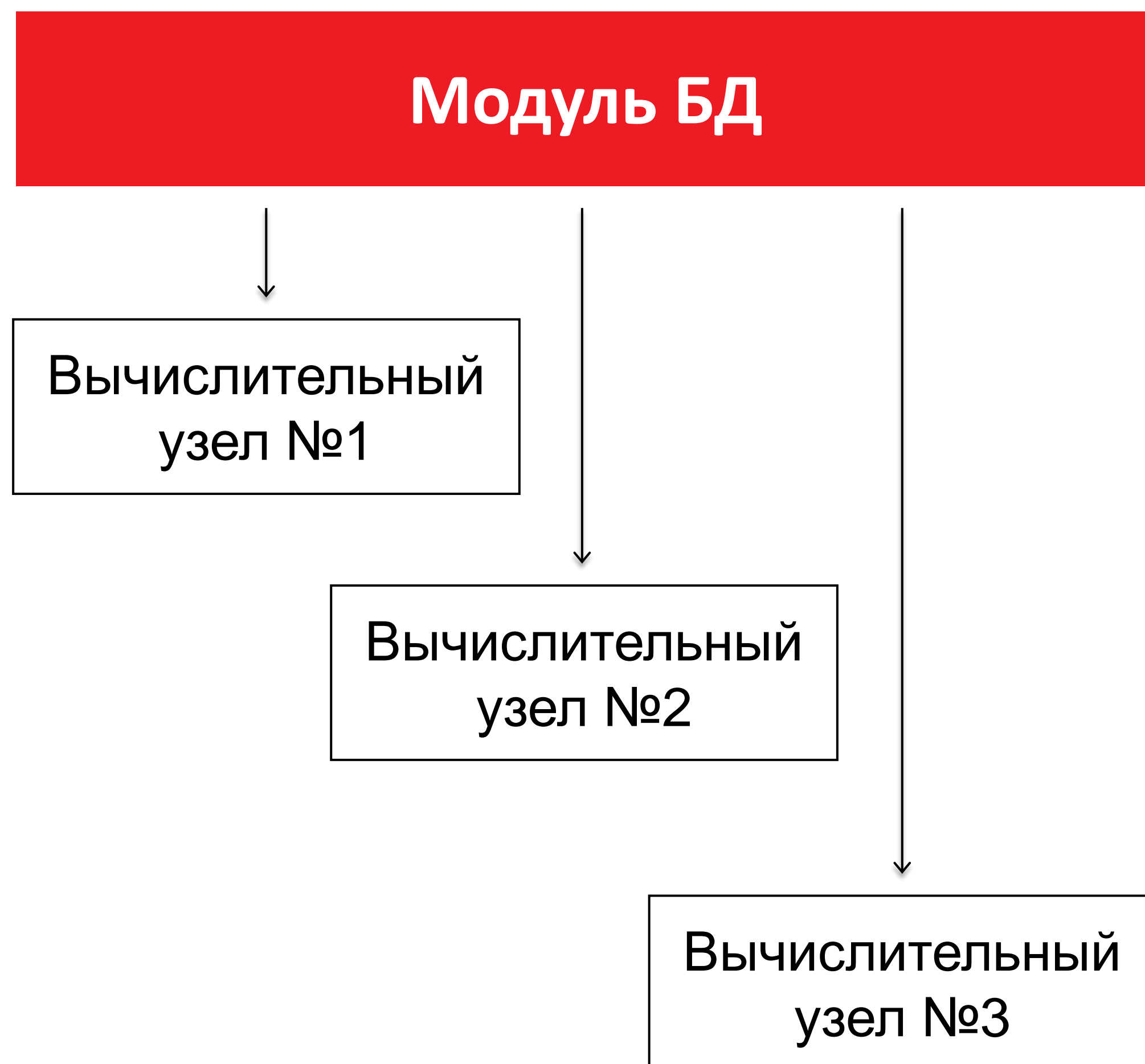
Сетевой компонент:

- Сетевой узел доступа (36 портов 10|25 ГБ/с + 4 порта 100 ГБ/с)
- Сетевой узел интерконнекта (28 портов 100 ГБ/с)

Компонент управления:

- Служебный узел
- Сетевой узел управления

Снаружи и внутри



**3 вычислительных узла,
в каждом:**

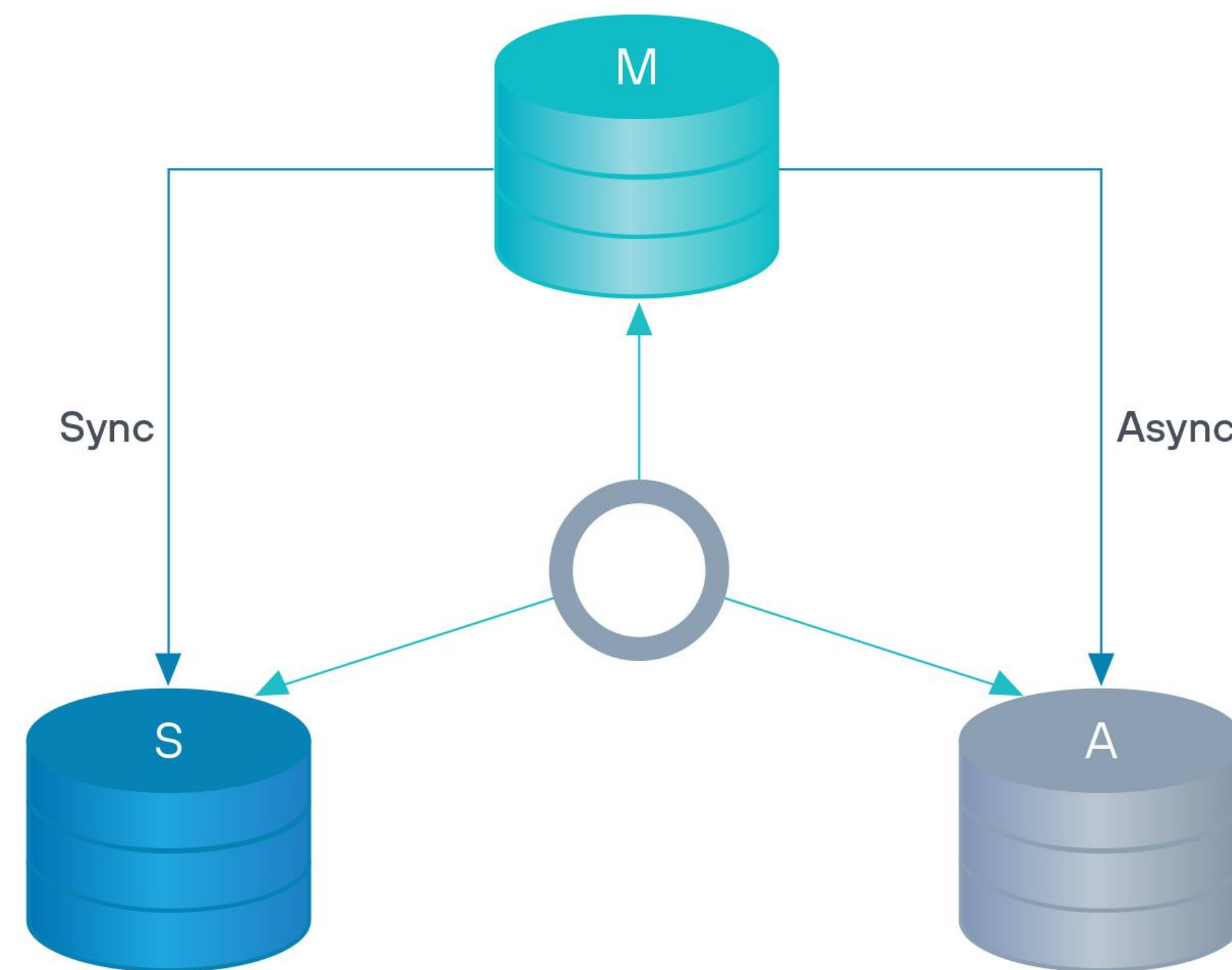
- 48 ядер 3.0 ГГц
- 768 ГБ оперативной памяти
- 10-20 ТБ дискового пространства для БД

Снаружи и внутри



Кластерное ПО

- Автоматическая обработка однократных отказов
- Сохранение точки подключения
- Сохранность данных и сервиса СУБД при отказе одного узла
- Автоматическая блокировка доступа в случае отказа двух узлов для гарантированной сохранности данных



Сетевое взаимодействие



Разделение сетей:

Сети публичного сервиса

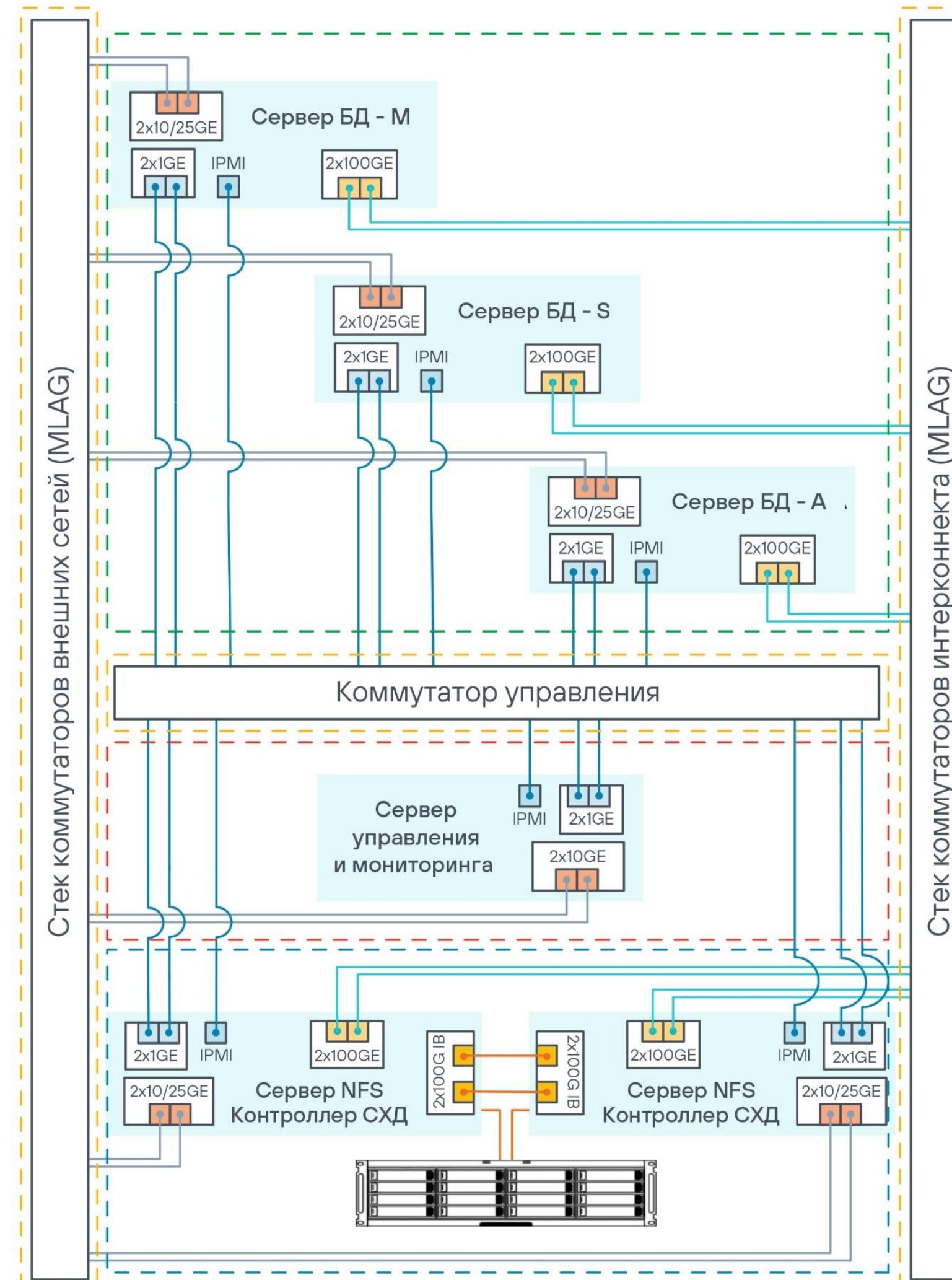
- Сеть доступа (public - default)
- Сеть внешнего бэкапа (backup)
- Сеть георепликации (dci)

Сети управления

- Только на служебном сервере
- IPMI, MGMT (default)

Внутренние сети

- Сеть интерконнекта (internal)
- Сеть установки (PXE)
- Сеть IPMI_C (web-proxy)
- Сеть Ring (web-proxy)



Скала МБД.П вместо Oracle Exadata

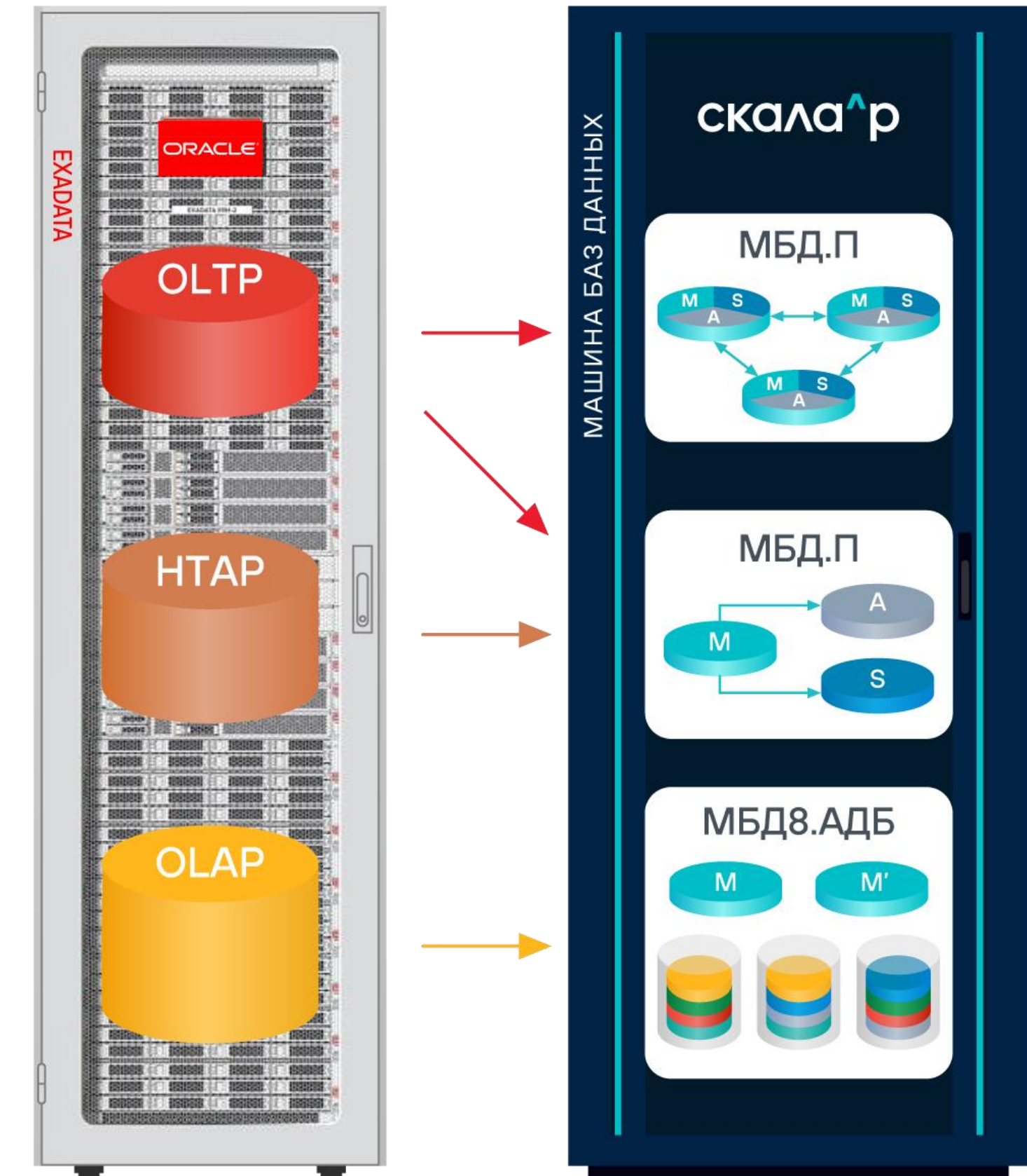


Oracle Exadata консолидирует 3 типа нагрузки:

- Транзакционная (OLTP)
- Гибридная (HTAP)
- Аналитическая (OLAP)

Транзакционная и гибридная типы нагрузок мигрируют на Скала^р МБД.П с максимально возможным уровнем производительности, доступности, сохранности данных:

- OLTP – до 3 сервисов СУБД на кластер
- HTAP – 1 сервис СУБД, использование синхронной реплики



Типовые комплекты поставки



Модель	М-1	М-2	М-3
Параметр			
Кол-во вычислительных модулей	1	2	3
Кол-во узлов БД	3	2x3	3x3
Кол-во БД	До 3	До 6	До 9
Общий объем БД, ТБ	До 20	До 2x20	До 3x20
Объем хранения СРК, ТБ	До 130	До 260	До 260



Скорость дисковой подсистемы



Программный рейд

- Производительнее аппаратного RAID-контроллера
- Минимальное использование RAM (требуется менее 4GB RAM)
- Управление процессорными потоками
- Минимальная просадка производительности в режиме восстановления
- Мониторинг износа накопителя
- Автозамена накопителя
- Повышение уровня RAID

```
fio -ioengine=libaio -direct=1 -bs=4k -  
iodepth=1 -fsync=1 -rw=randwrite
```

ERA Raidix RAID10:

- write: IOPS=20.7k, BW=80.9MiB/s (84.8MB/s);
- lat (usec): min=34, max=756, avg=44.6

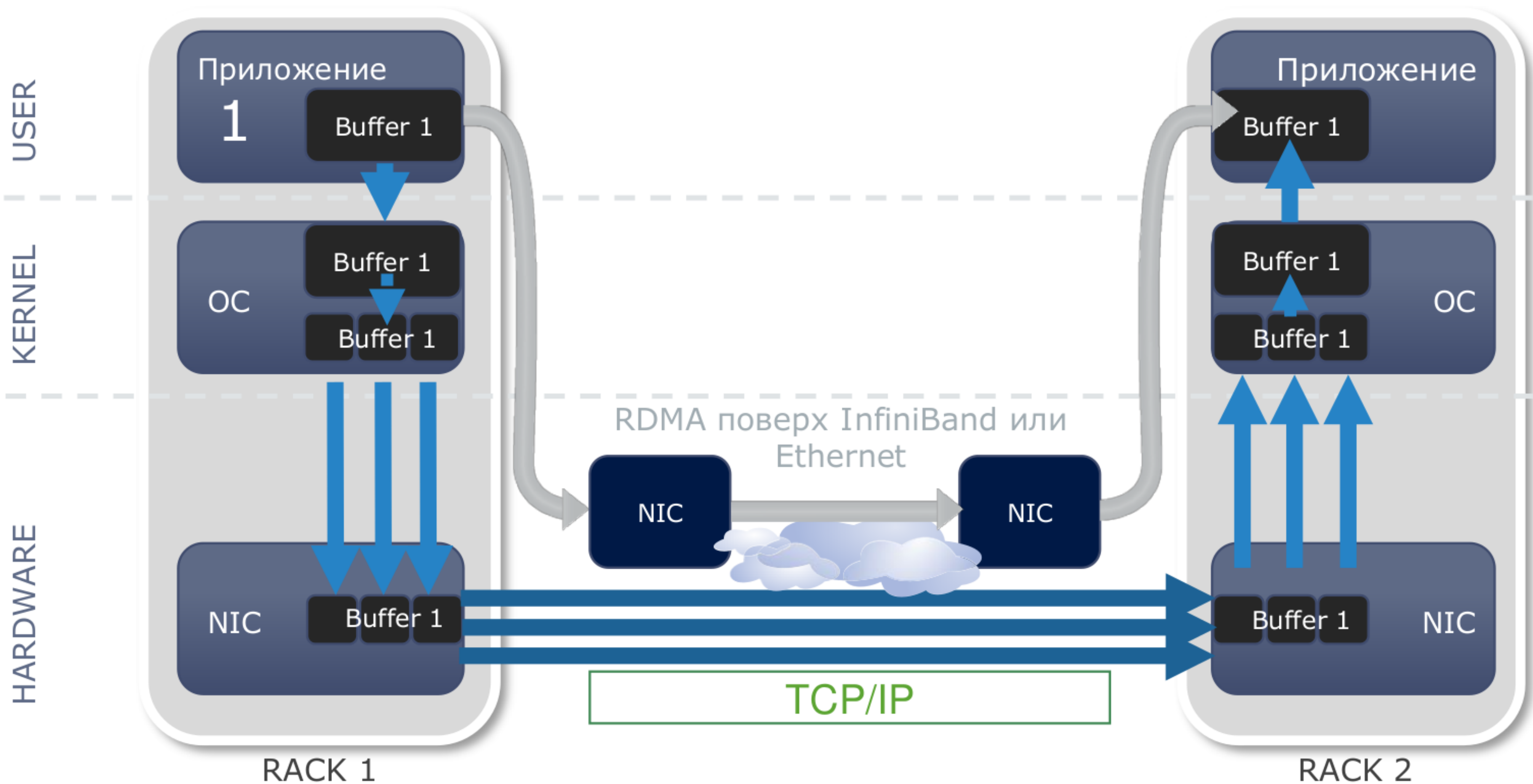
mdadm RAID10:

- write: IOPS=16.8k, BW=65.5MiB/s (68.7MB/s);
- lat (usec): min=35, max=671, avg=55.52

LSI MegaRAID RAID10:

- write: IOPS=19.1k, BW=73.9MiB/s (77.4MB/s);
- lat (usec): min=41, max=528, avg=48.27

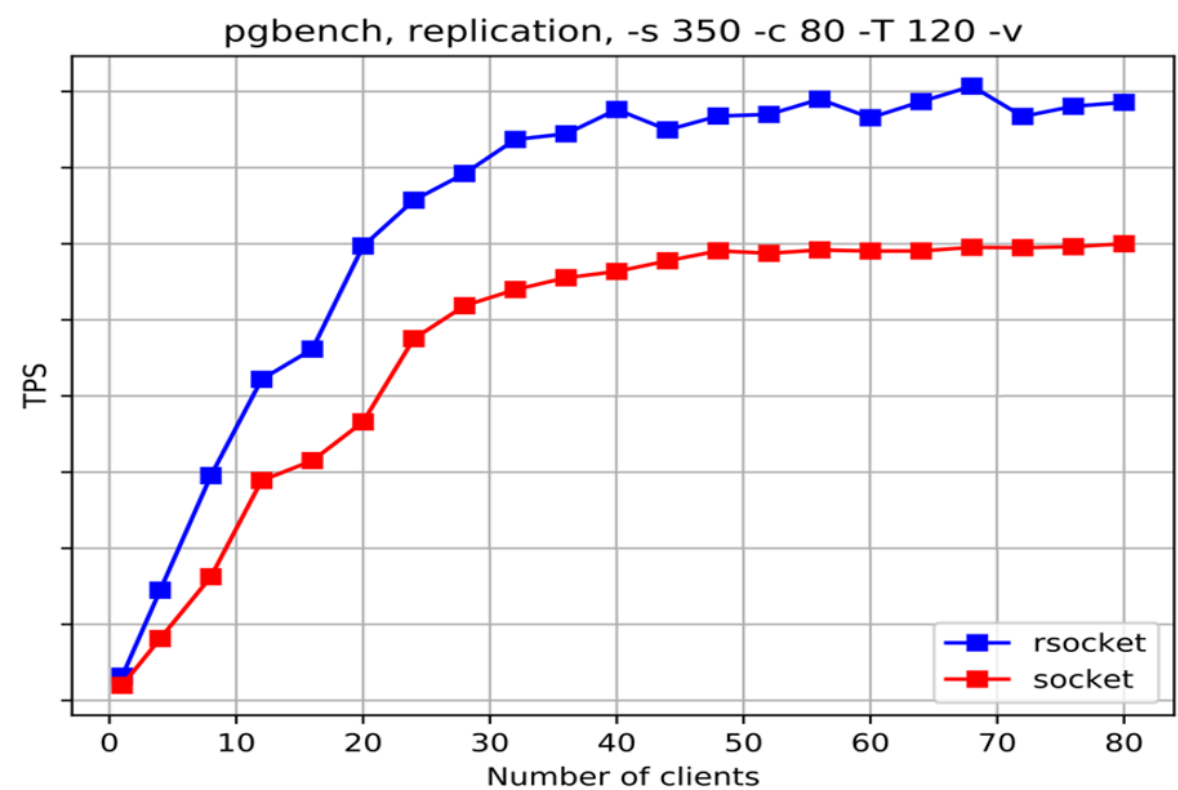
Сети. Использование RDMA



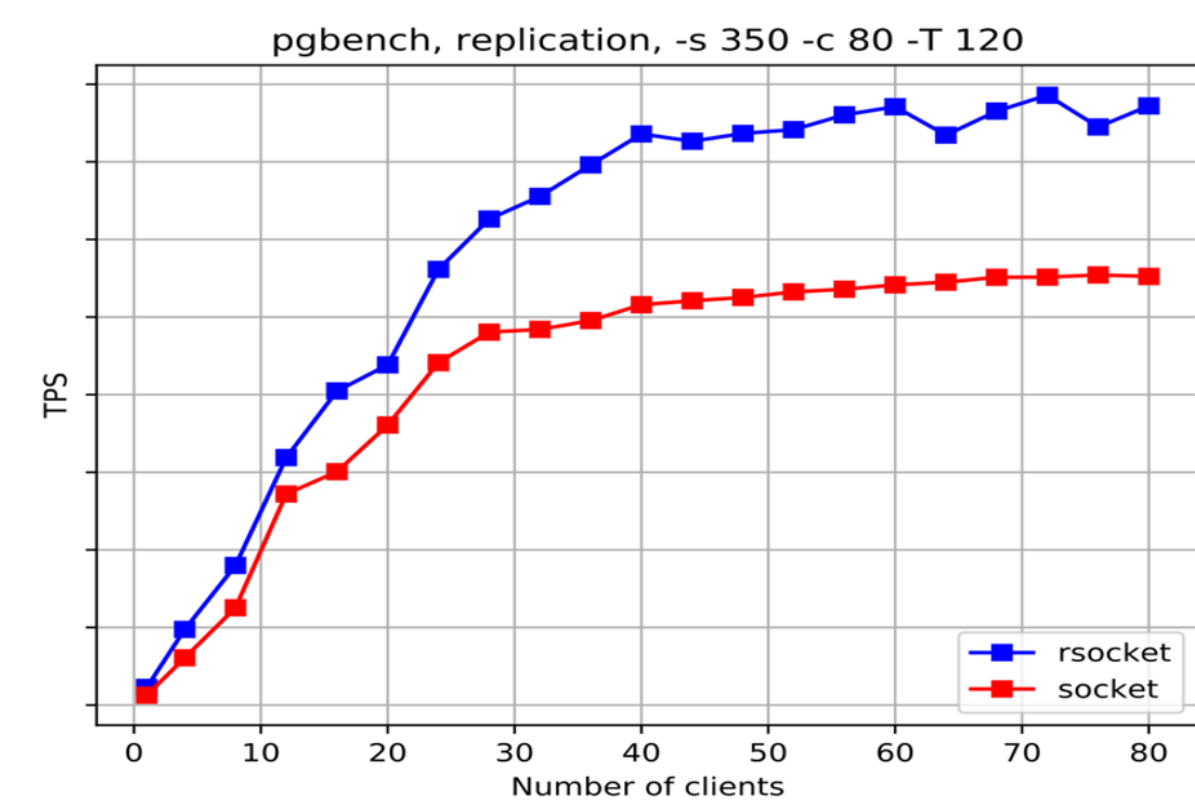
Сети. Использование RDMA



Read-write with synchronous replication, TPS



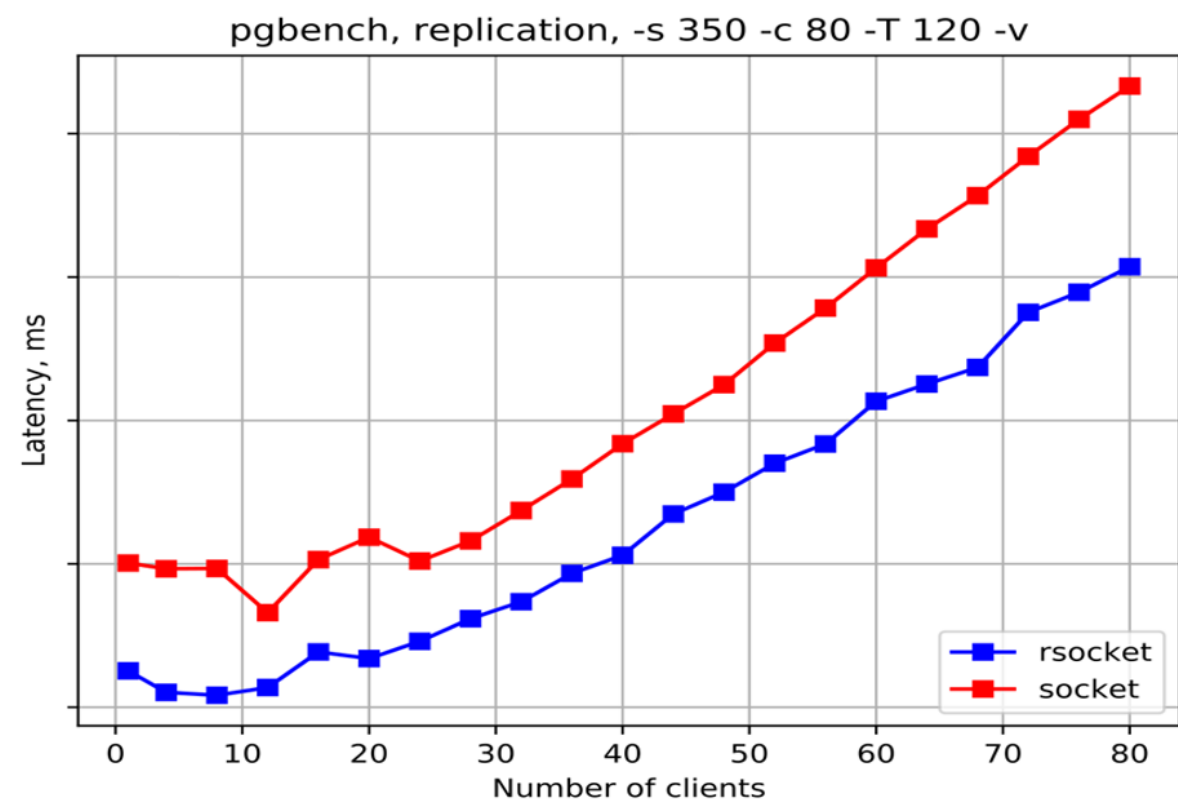
Read-write with remote apply replication, TPS



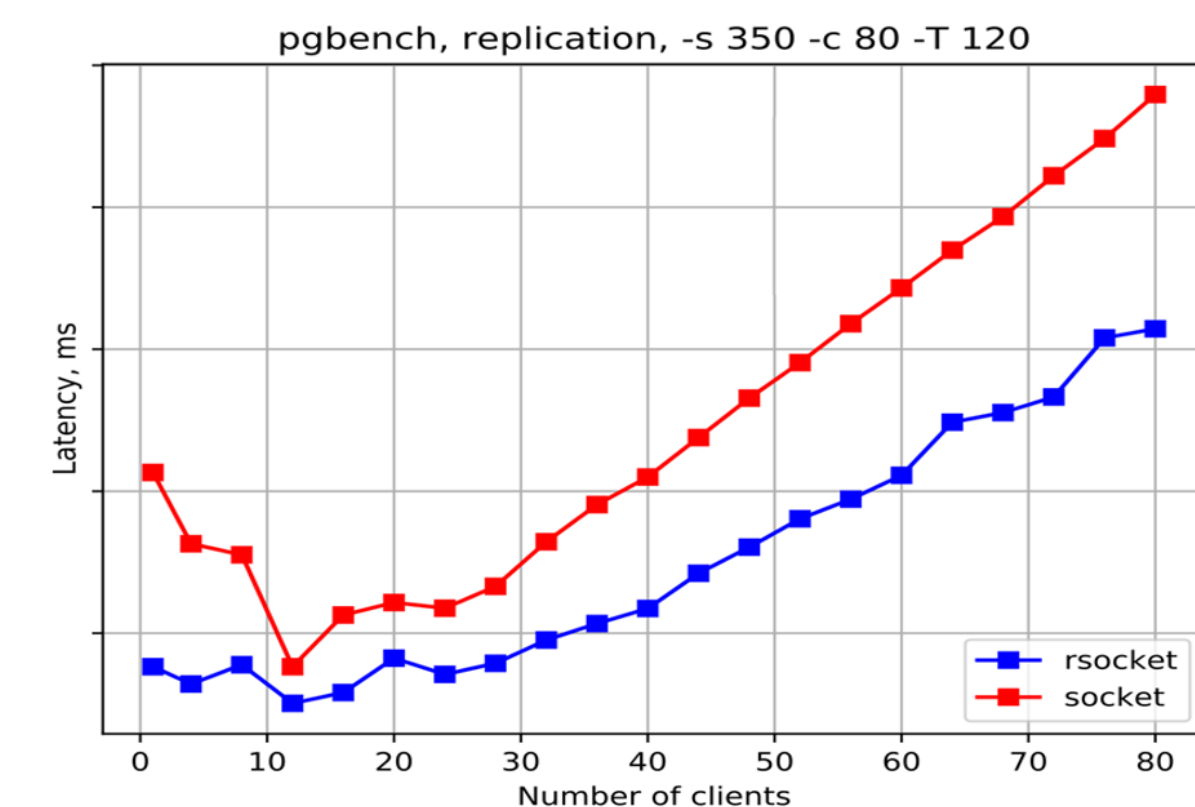
Знаешь Скала^р?
Отметься!



Read-write with synchronous replication, latency



Read-write with remote apply replication, latency



Сети. Не все коммутаторы одинаково полезны



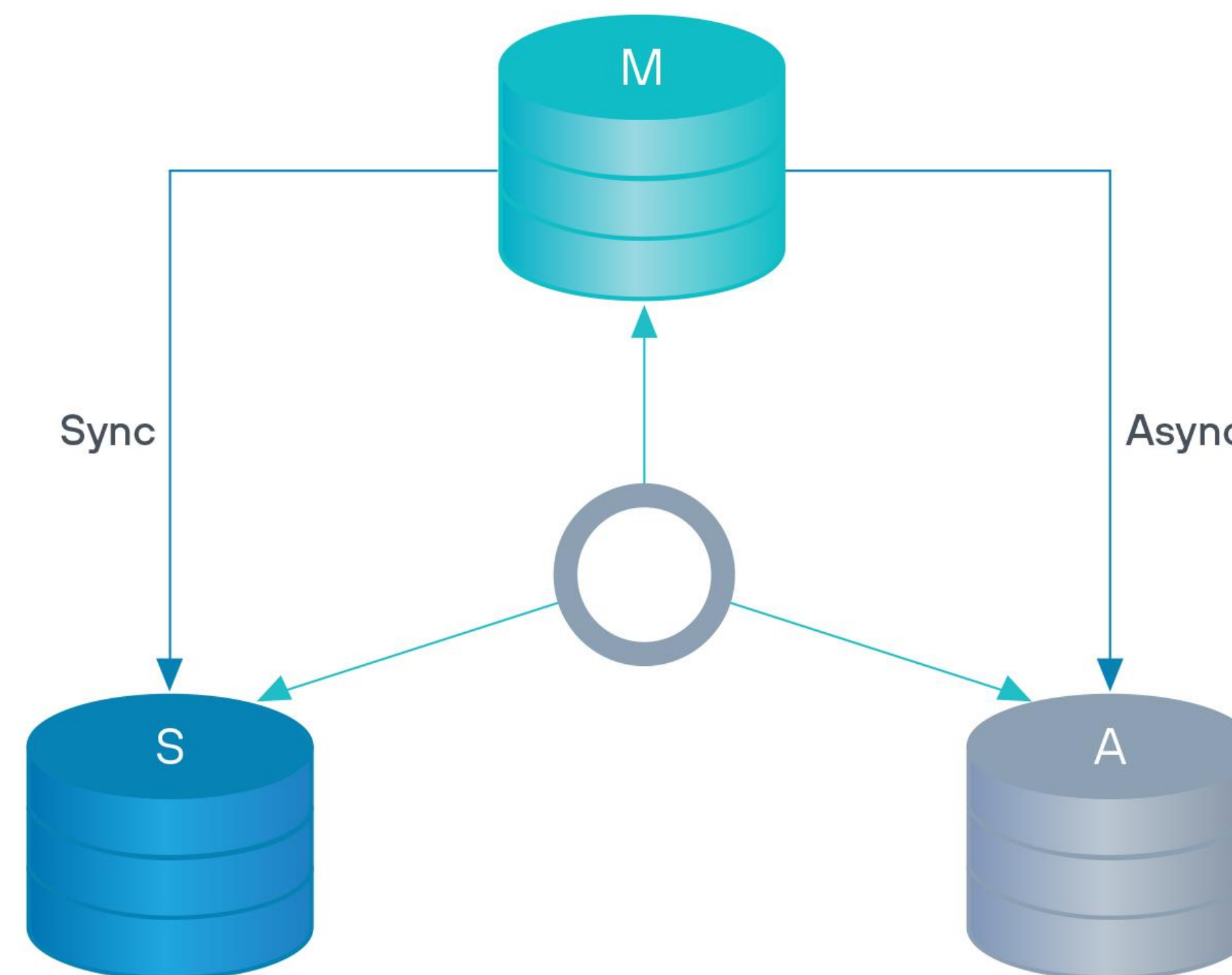
Производитель	Модель	Время сходимости	Пропускная способность	VLAN	RSTP	RDMA	Порты	Производство	Цена
Vendor#1	Model#1	100	92	4000	On	+/-	32	РФ	0.8X
Vendor#2	Model#2	2000	96	4096	Off	-	32	РФ	2.7X
Vendor#3	Model#3	3000	100	4096	?	+	32	Китай	2.5X
Vendor#4	Model#4	250	100	4096	On	+	32	Китай	1.5X
Mellanox	SN2100	40	100	4096	On	+	16	США	X
Vendor#5	Model#5	-	100	4000	On	-	32	РФ	???

Отказоустойчивость



Отказоустойчивость на всех уровнях

- Надежные комплектующие
- Резервирование значимых компонентов на аппаратном уровне
- Отказоустойчивая архитектура верхнего уровня
 - Оперативное восстановление при сбоях
 - Автоматическая обработка однократных отказов
- Сохранение точки подключения
- Сохранность данных и сервиса СУБД при отказе одного узла
- Автоматическая блокировка доступа в случае отказа двух узлов для гарантированной сохранности данных

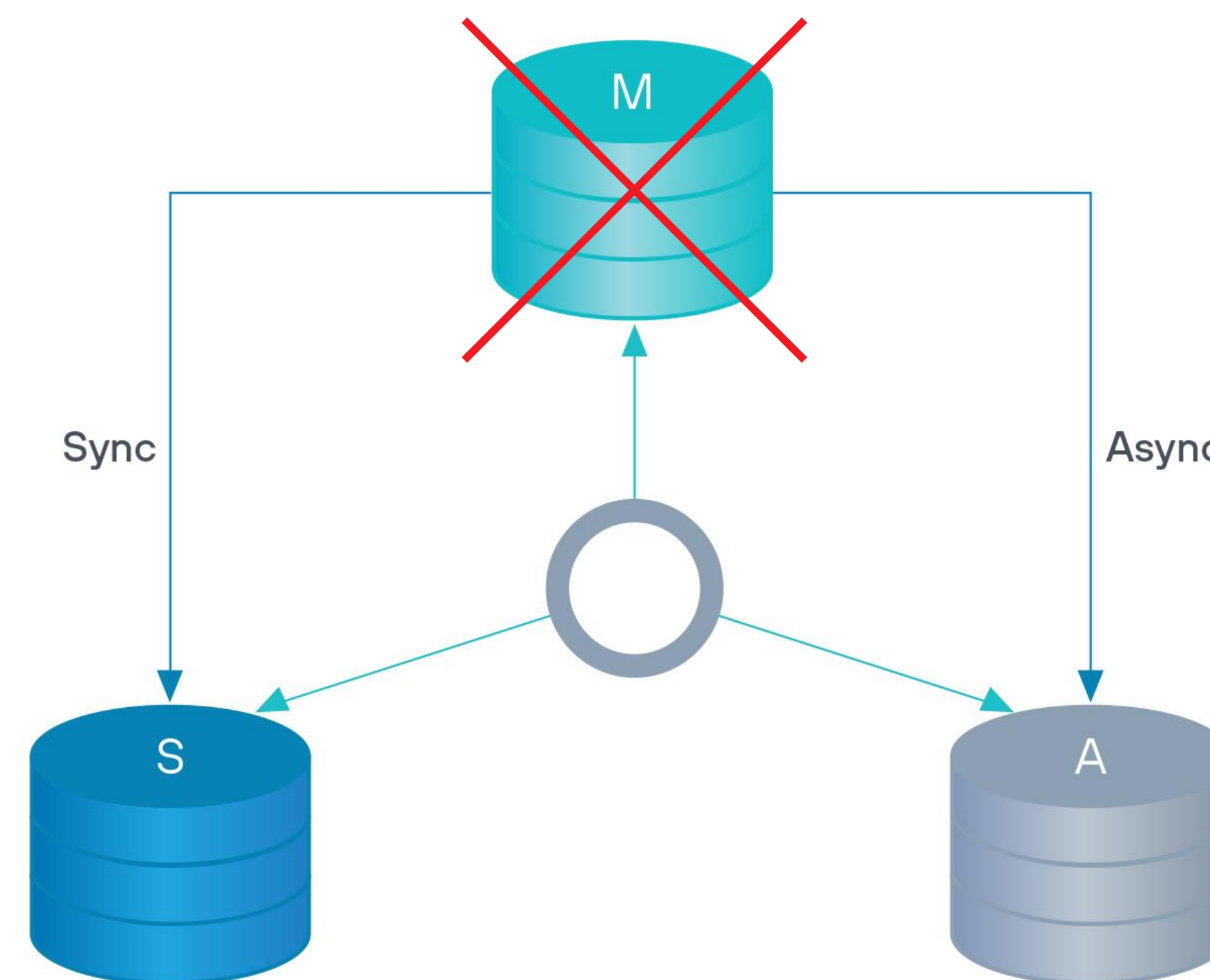


Отказоустойчивость



Отказоустойчивость на всех уровнях

- Отказоустойчивая архитектура верхнего уровня
 - Оперативное восстановление при сбоях
 - Автоматическая обработка однократных отказов
- Сохранение точки подключения

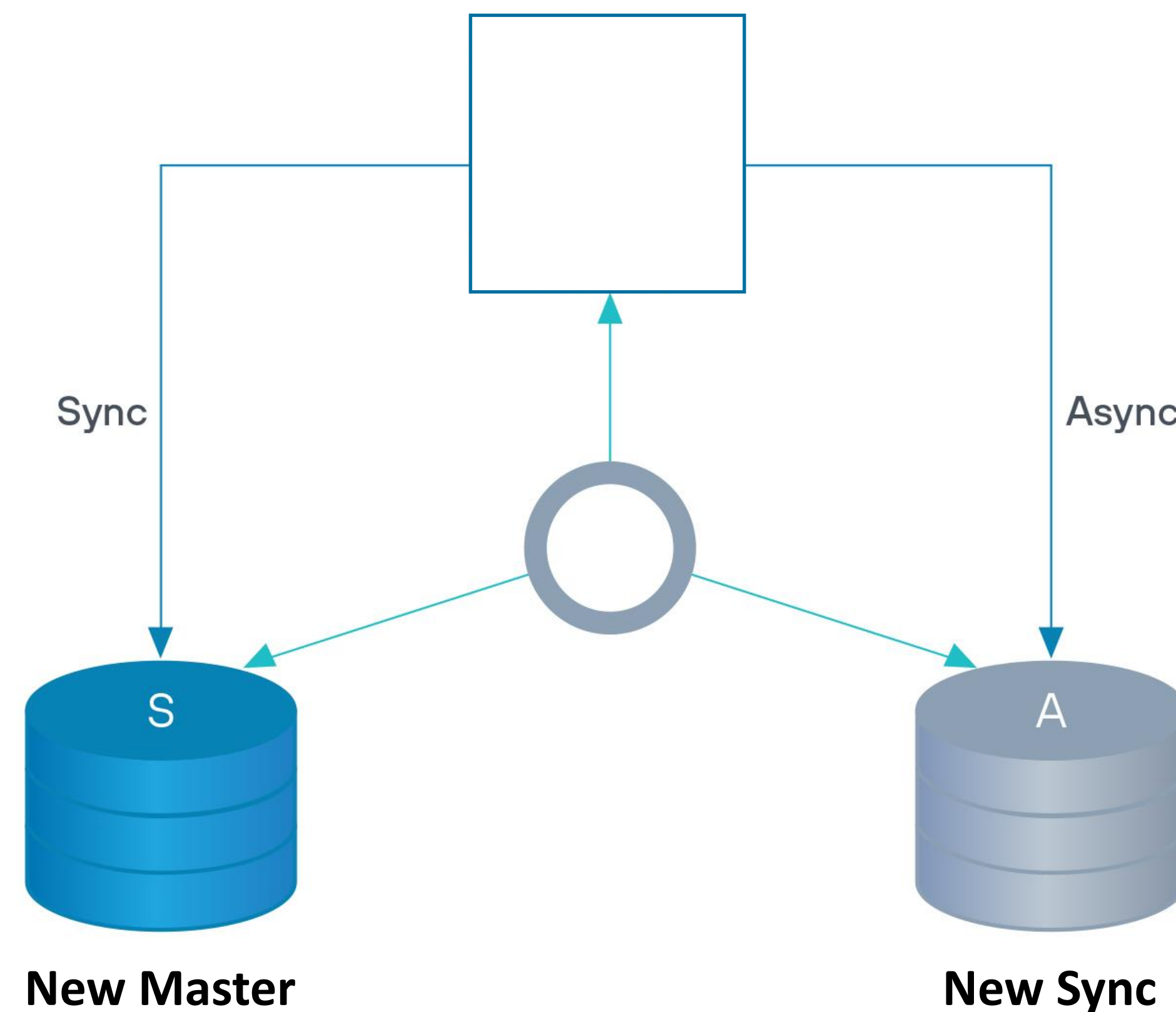


Отказоустойчивость

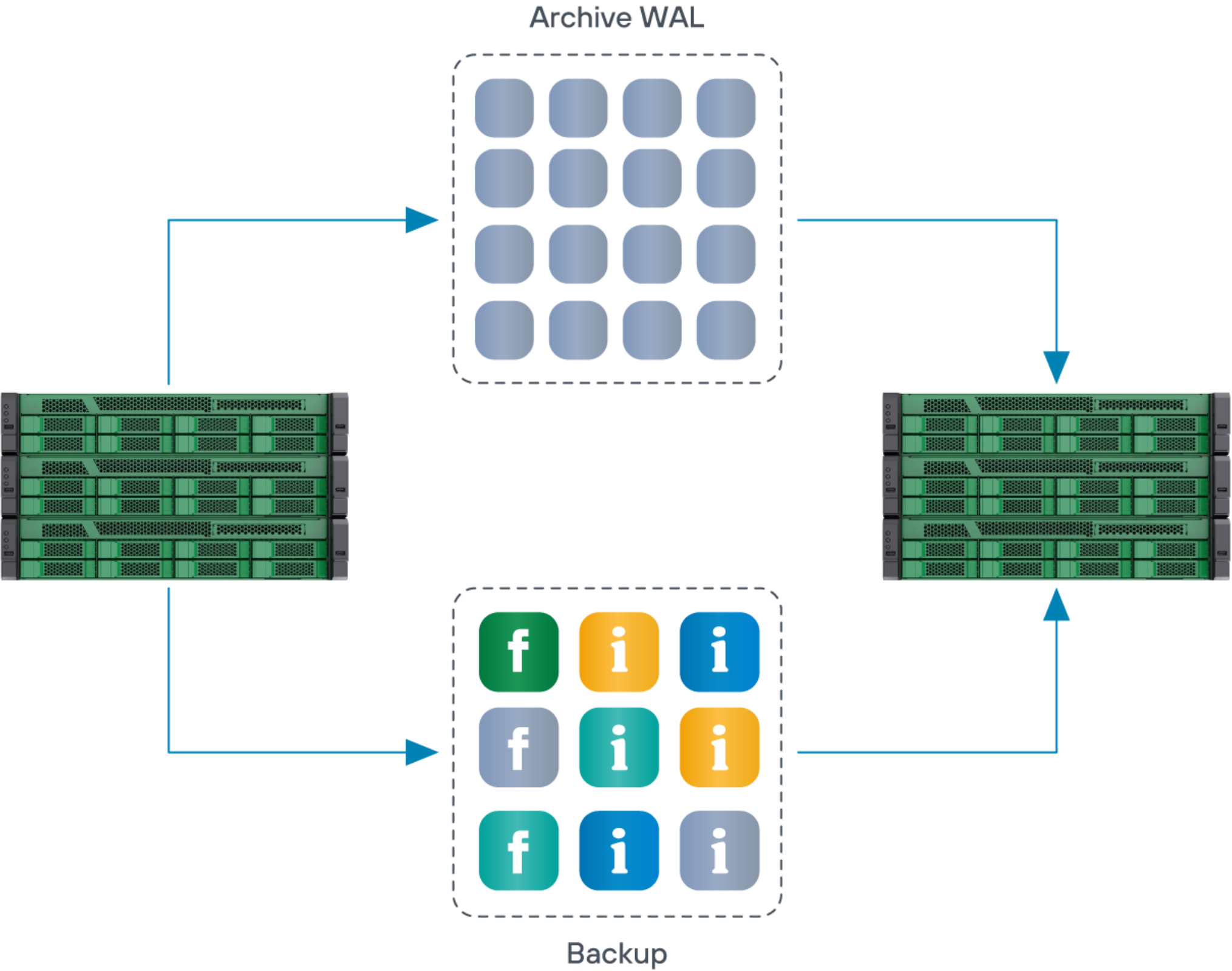
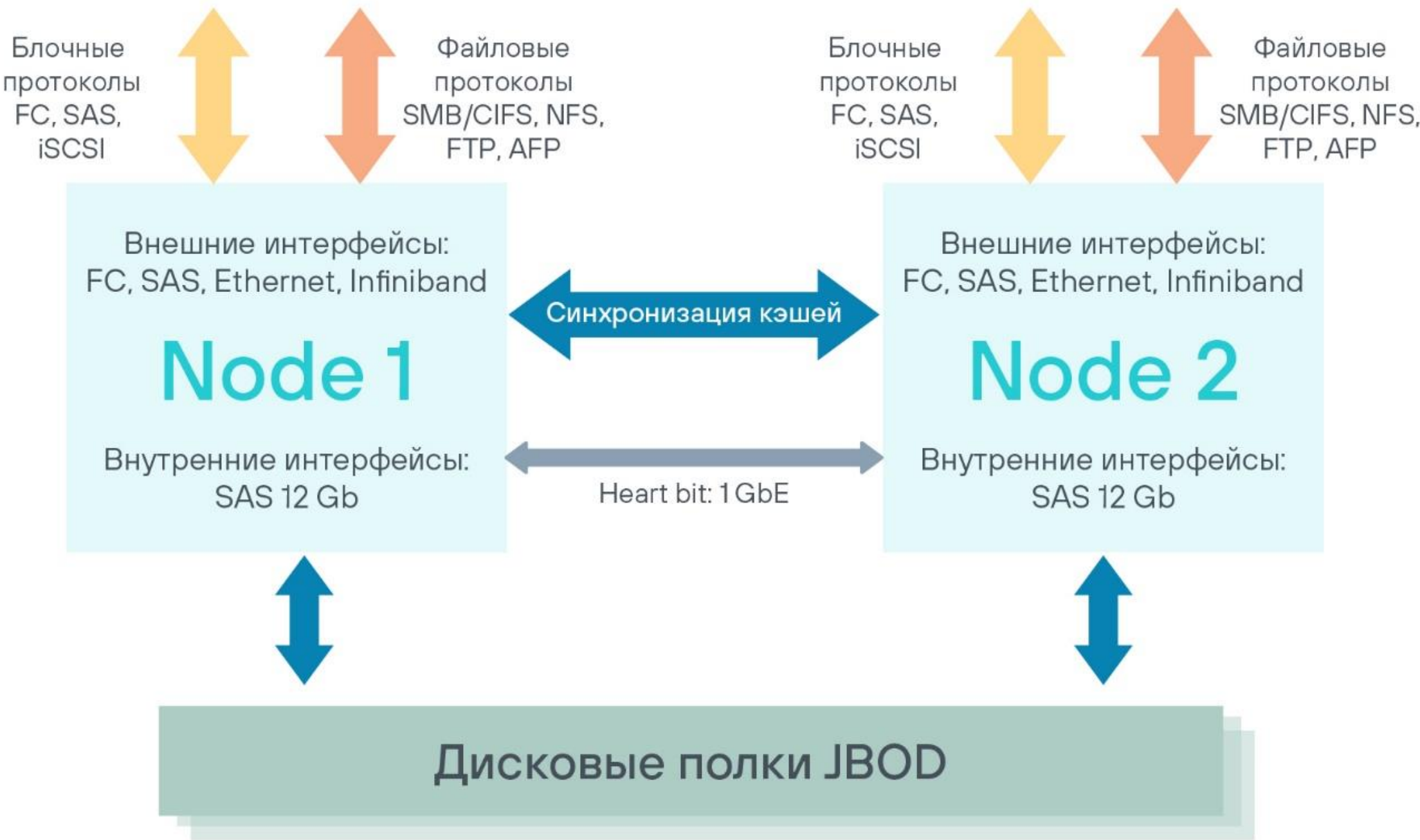


Отказоустойчивость на всех уровнях

- Сохранение точки подключения
- Сохранность данных и сервиса СУБД при отказе одного узла
- Автоматическая блокировка доступа в случае отказа двух узлов для гарантированной сохранности данных



Резервное копирование



Меряем производительность



Измеряем правильно!

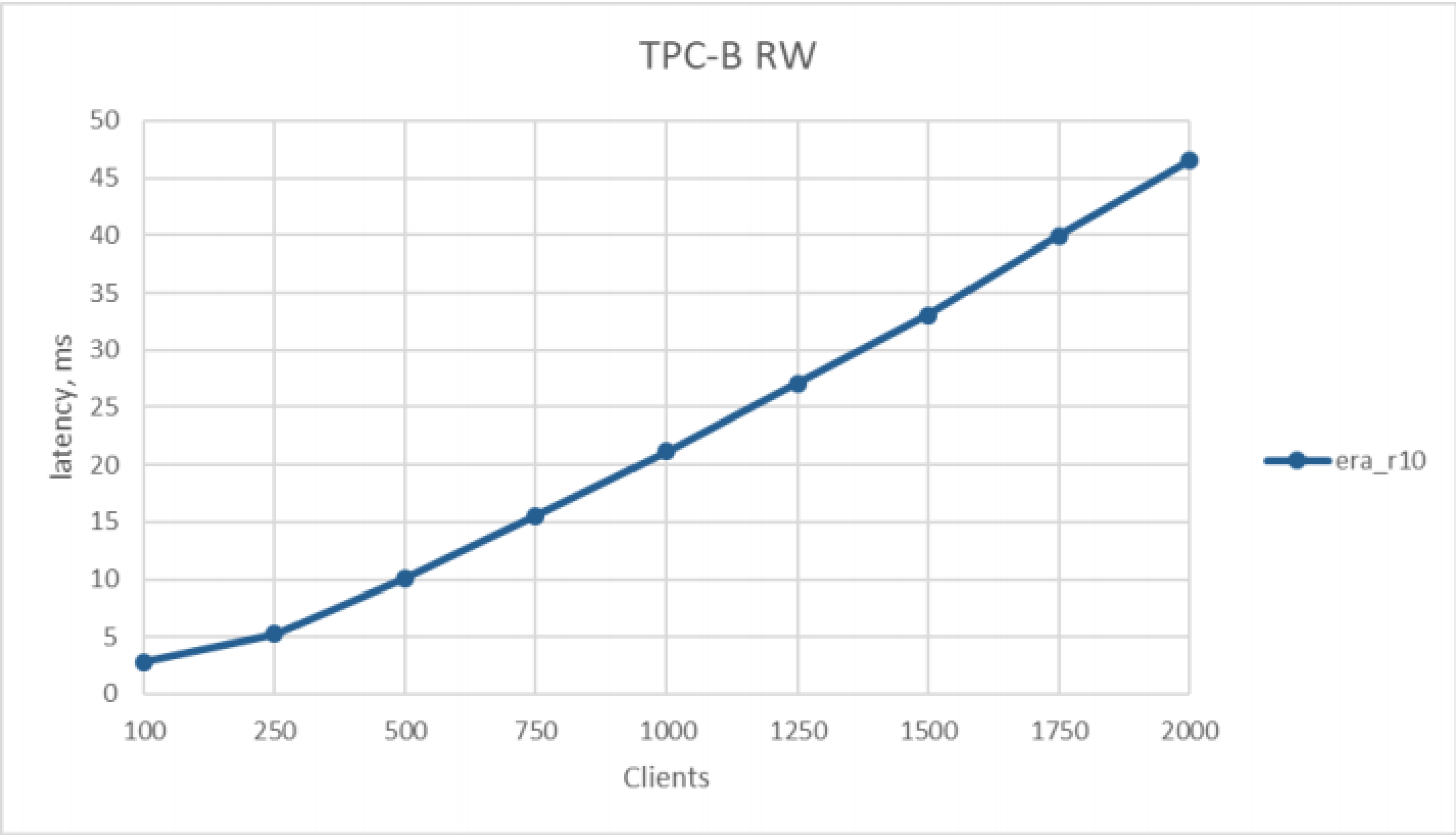
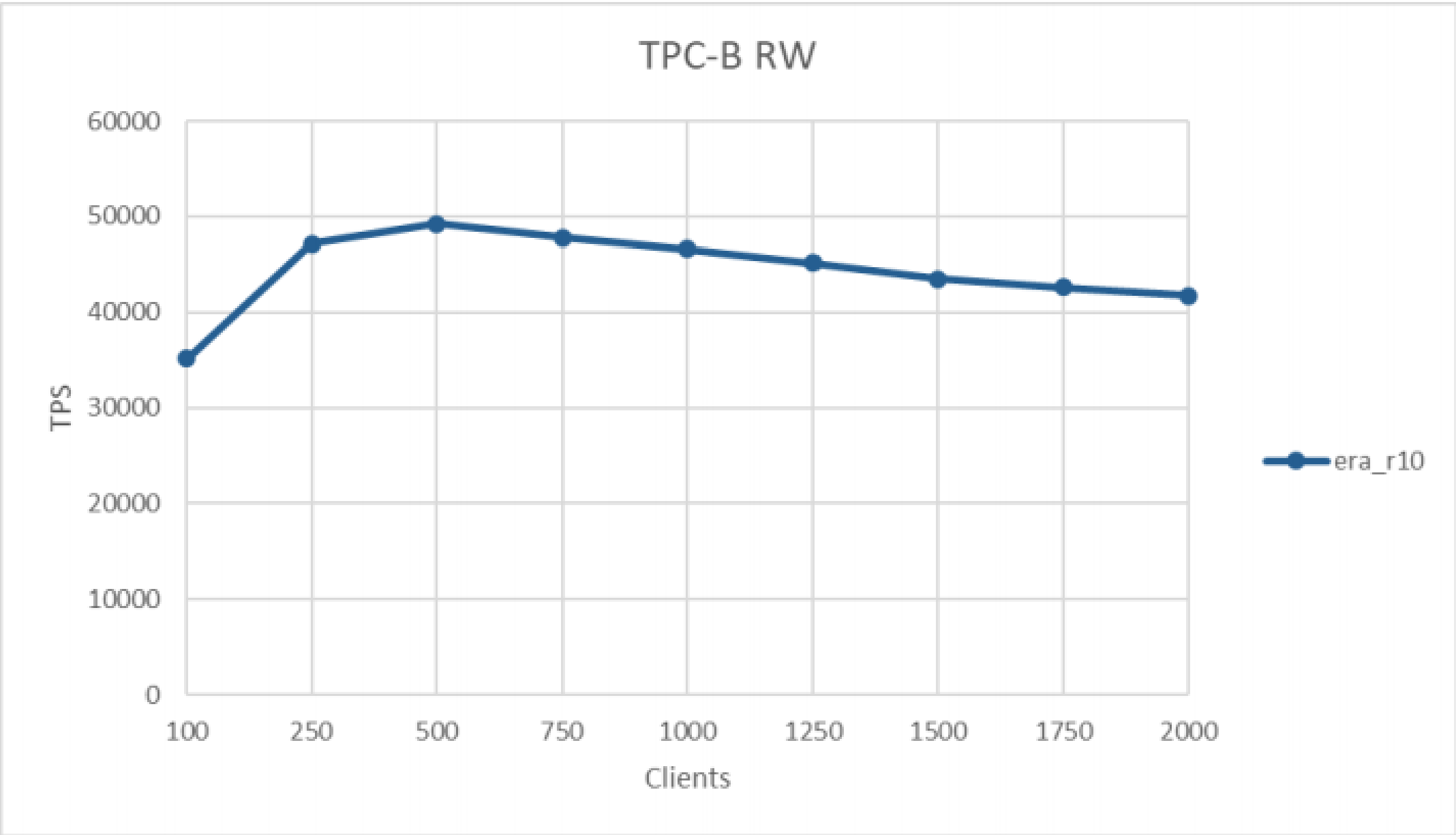
Андрей Николаенко, PG Day'17 –
«Об эталонном тестировании PostgreSQL»



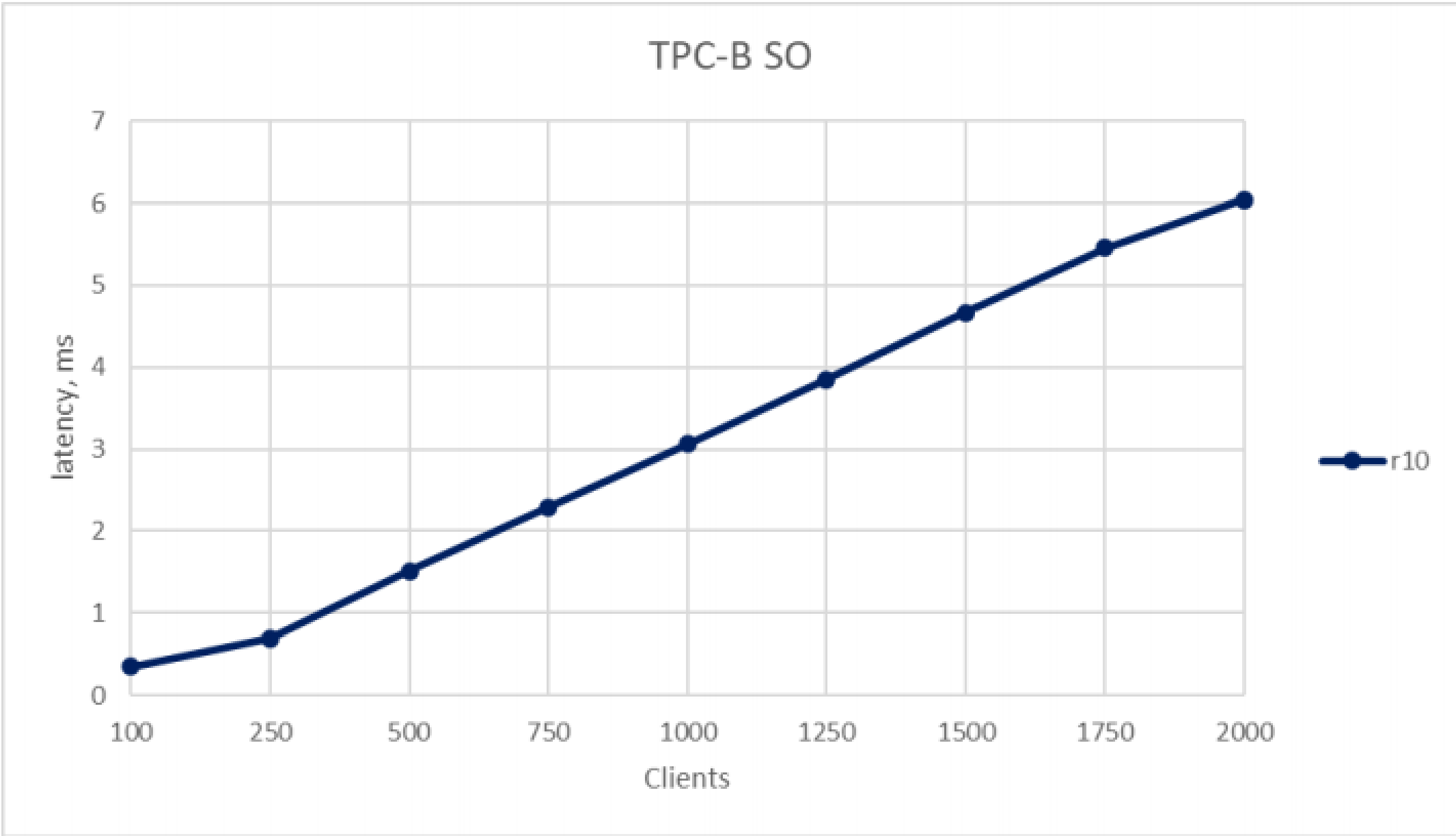
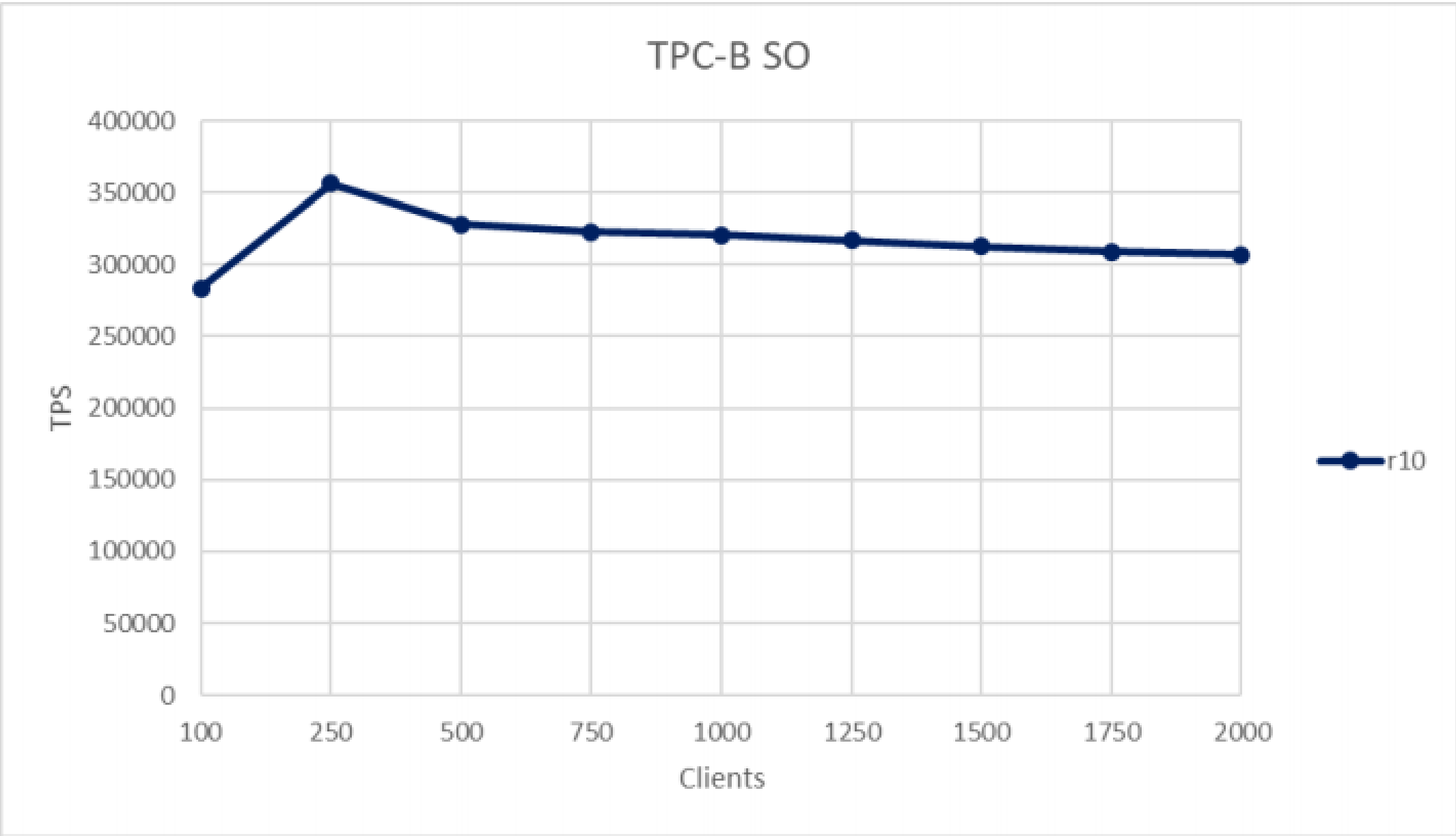
- pgbench, TPC-B-подобная нагрузка
- SF=70000
- RW-тест с БД объемом более 1 Тб
- SO-тест с тома объемом 12 Тб

```
1.BEGIN;  
2.UPDATE pgbench_accounts SET abalance = abalance + :delta WHERE aid = :aid;  
3.SELECT abalance FROM pgbench_accounts WHERE aid = :aid;  
4.UPDATE pgbench_tellers SET tbalance = tbalance + :delta WHERE tid = :tid;  
5.UPDATE pgbench_branches SET bbalance = bbalance + :delta WHERE bid = :bid;  
6.INSERT INTO pgbench_history (tid, bid, aid, delta, mtime) VALUES (:tid, :bid, :aid, :delta, CURRENT_TIMESTAMP);  
7.END;
```


Меряем производительность



Меряем производительность



- Вечерний пятничный полный бэкап ~11 ТБ (сжатие) за 16,5 часов (режим ARCHIVE) + 500 ГБ архивных журналов
- Утренний инкрементальный бэкап в понедельник ~200–400 ГБ за 30–50 минут (режим PAGE) + 10–20 ГБ архивных журналов
- Инкрементальный бэкап каждые 2 часа в рабочие дни ~10–30 ГБ за ~3–5 минут (режим PAGE) + 5–15 ГБ архивных журналов



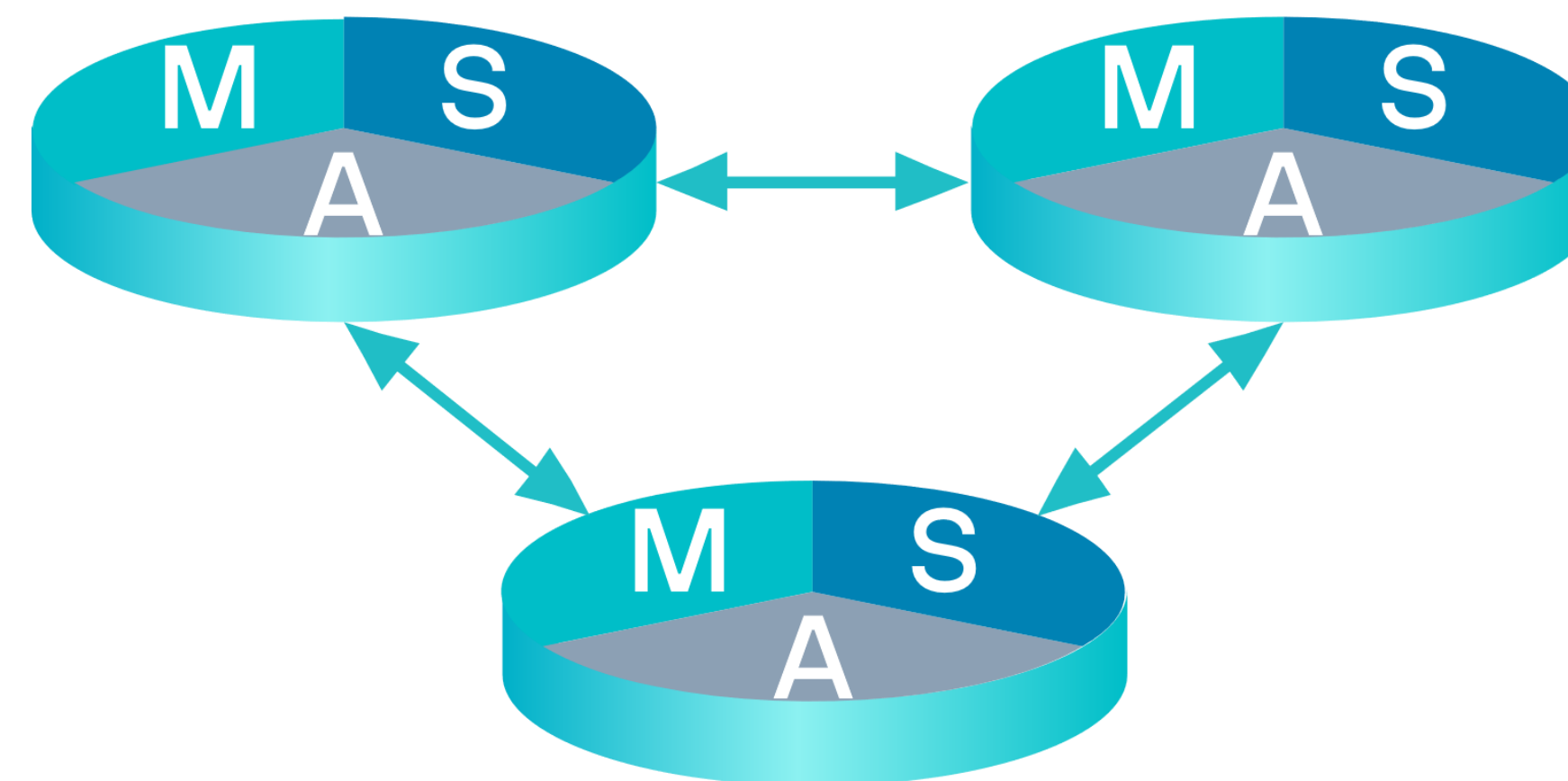
Случай из жизни №2. «Шахматка»



Решение – Скала^р МБД.П:

- 1 аппаратный кластер (3 узла)
- 3 экземпляра БД
- СУБД Postgres Pro Enterprise
- Объем БД № 1 ~18 ТБ
- Объем БД № 2 ~17,5 ТБ
- Объем БД № 3 ~15 ТБ

МБД.П

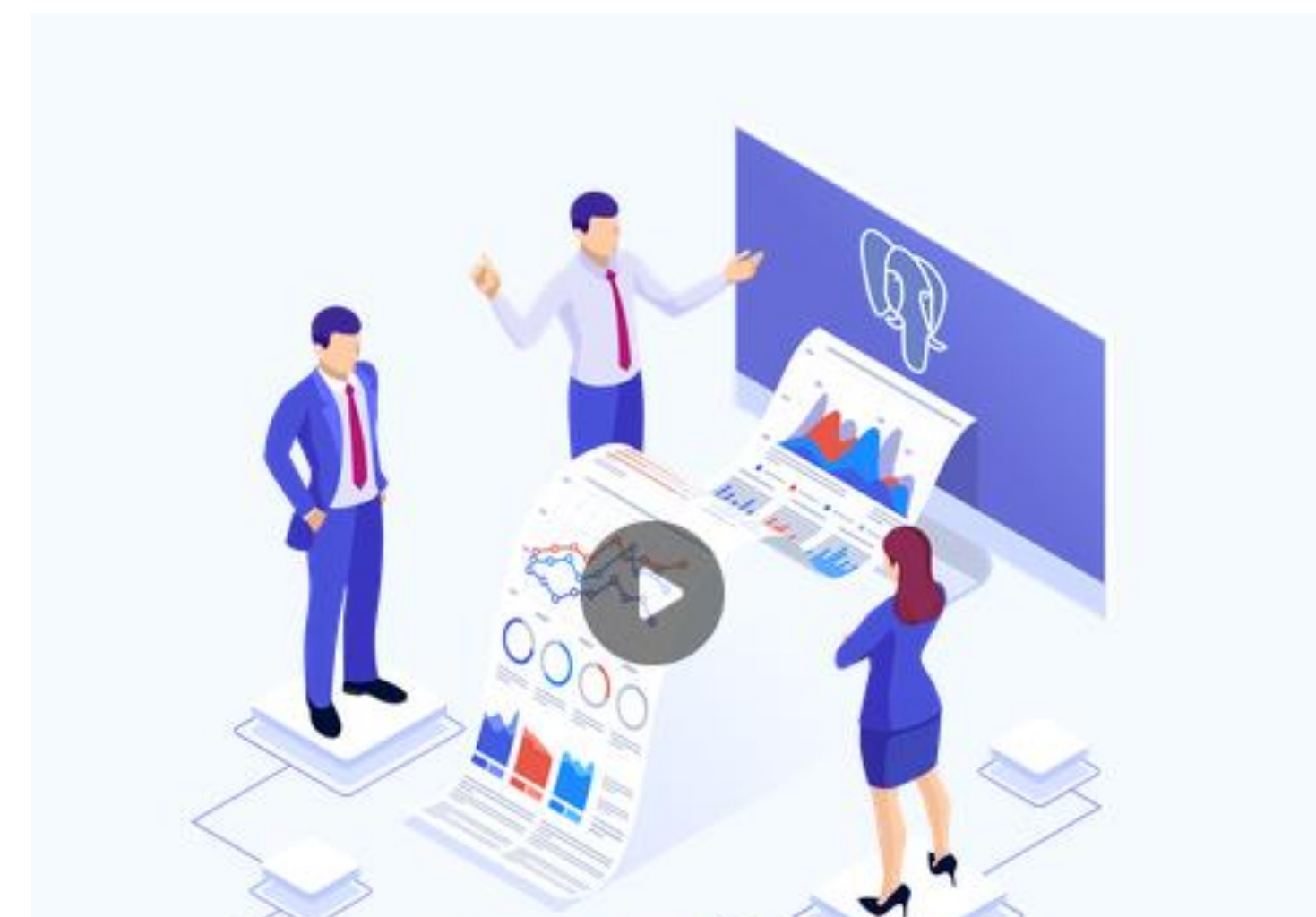


Система управления Машиной БД



Что есть:

- **SeveralNines Cluster Control**
<https://severalnines.com/product/clustercontrol>
- **ScaleGrid**
<https://scalegrid.io/postgresql/demo.html>
- **Awide**
<https://awide.io/>



PostgreSQL Pro

1. Create a PostgreSQL Cluster
2. Slow Query Analyzer
3. Backup and Restores
4. Dynamic Scaling

Create a PostgreSQL Dedicated Cluster

PostgreSQL Host
New Dedicated C

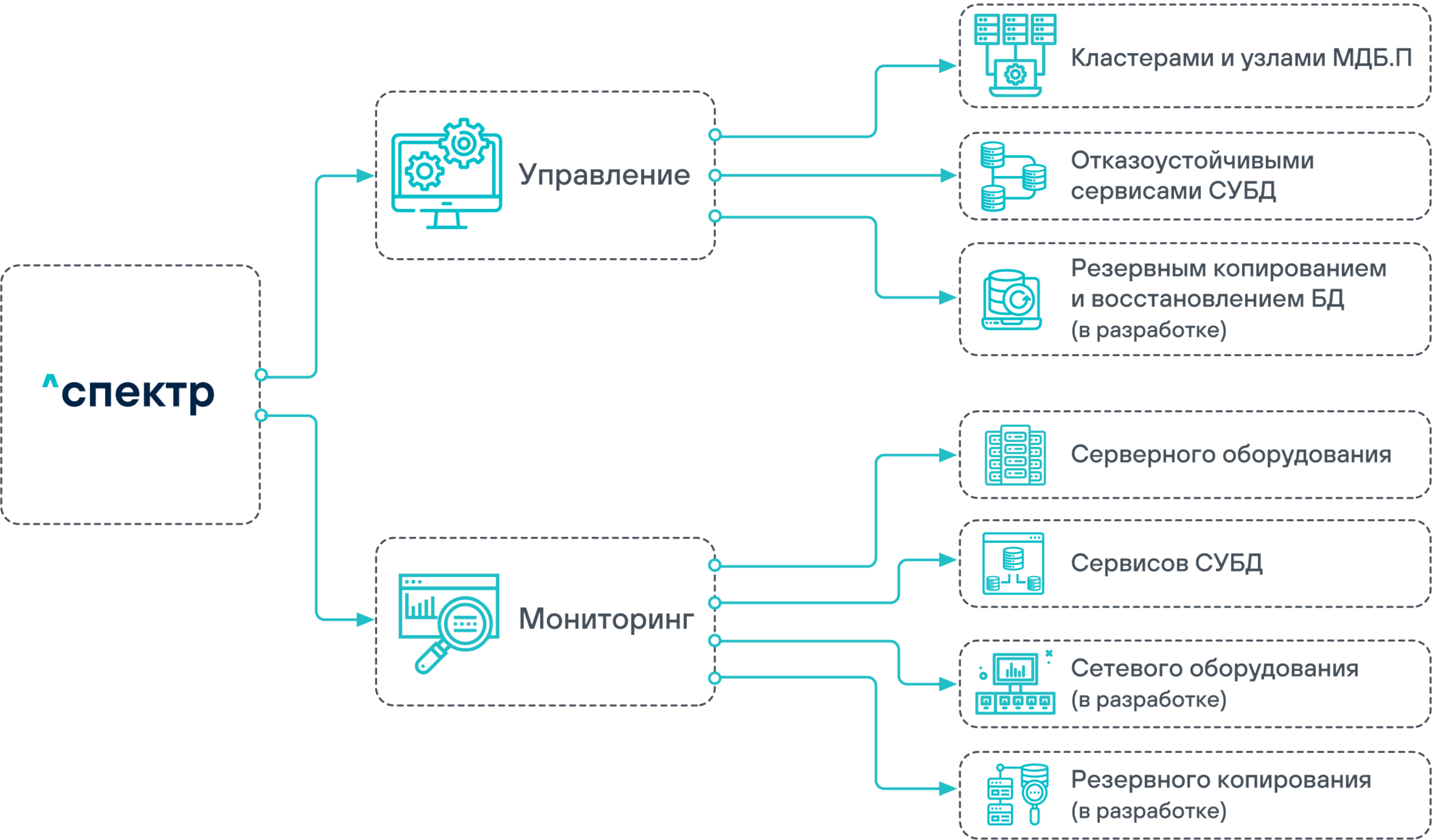
Hybrid multi-cloud database ops platform

ClusterControl empowers you to fully automate your open-source database operations on-premises, in the cloud, or both without giving up control.

Demo our sandbox!

Try free for 30 days

Система управления Машиной — ^Спектр!



Система управления Машиной — ^Спектр!



Снижает сложность операций с кластерным ПО
Использует проверенные шаблоны настроек

^спектр

ОБЪЕКТЫ УПРАВЛЕНИЯ

Кластеры

Узлы

Сервисы СУБД

ИСТОРИЯ

Операции

Кластеры

Поиск

+ Импортировать кластер

Кластер ↑↓	Статус ↑↓	Площадка ↑↓	↑↓	Узлы ↑↓	
alt82 ⓘ	ДОСТУПЕН	ЦОД 124	1 4 1 1	2 1	⋮
alt83 ⓘ	ДОСТУПЕН	ЦОД 123	1 5	1	
alt84 ⓘ	ДОСТУПЕН	ЦОД 123	1 3 1	3	⋮

Остановить кластер
Запустить кластер

^спектр

ОБЪЕКТЫ УПРАВЛЕНИЯ

Кластеры

Узлы

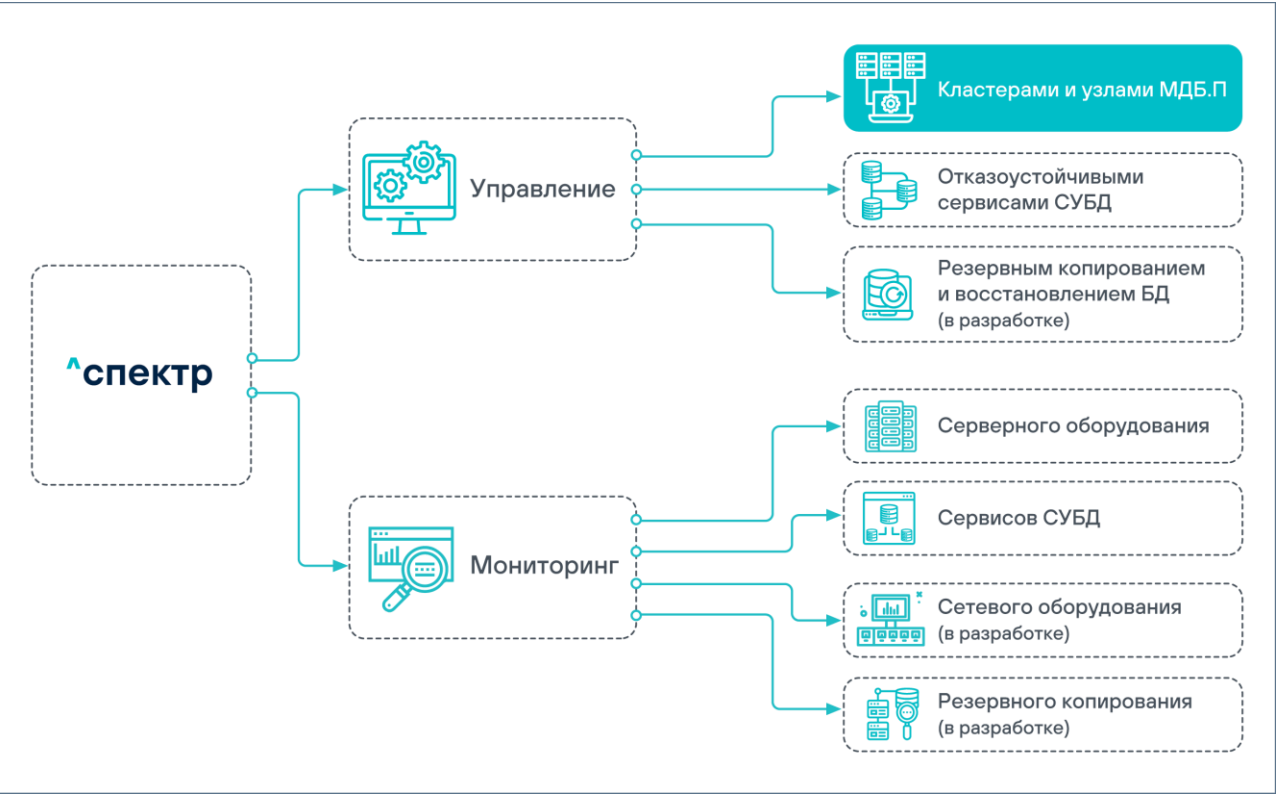
Сервисы СУБД

ИСТОРИЯ

Операции

Узлы

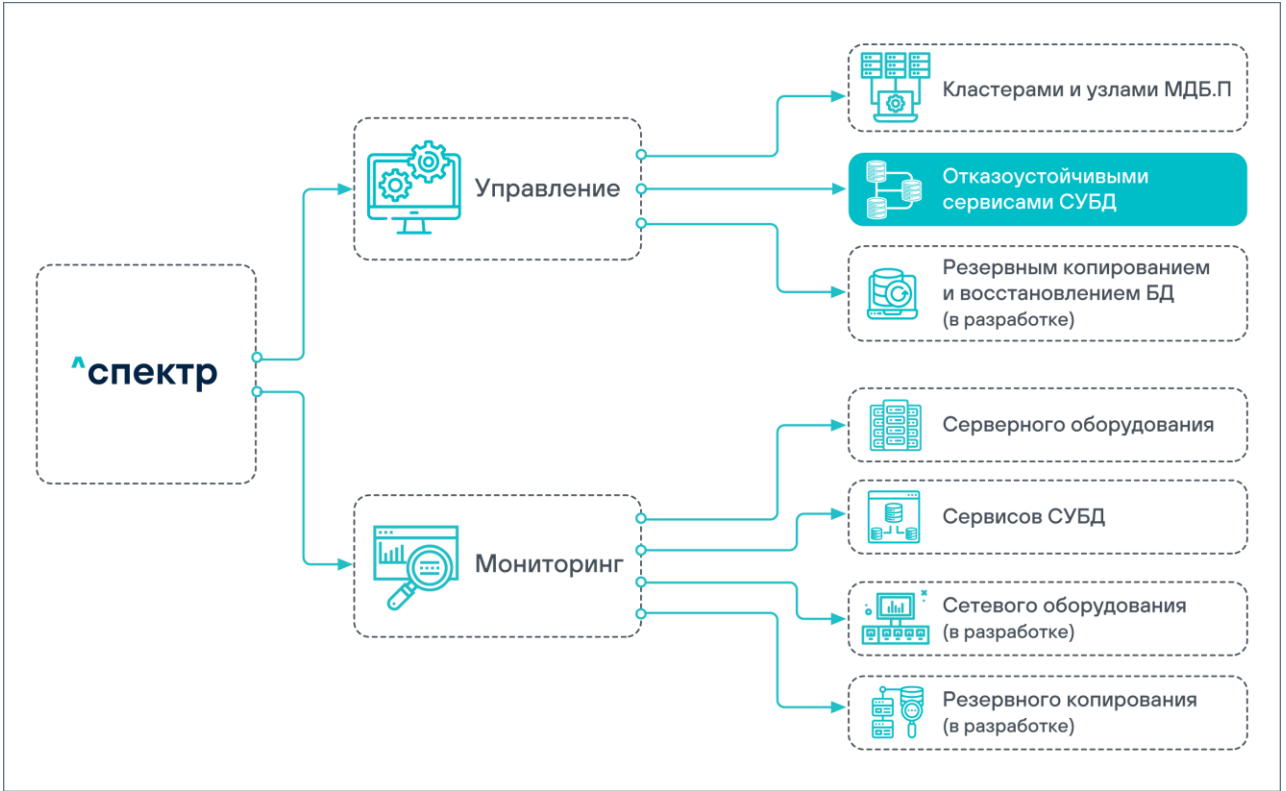
Узел ↑↓	IP / hostname ↑↓	Имя кластера ↑↓	Статус ↑↓	Площадка ↑↓	
alt82 node01 ⓘ	al8-node01	alt82	ДОСТУПЕН	ЦОД 123	⋮
alt82 node02 ⓘ	al8-node02	alt82	ДОСТУПЕН	ЦОД 123	⋮
alt82 node03 ⓘ	al8-node03	alt82	ДОСТУПЕН	ЦОД 123	⋮



Управление и настройка сервисов в СУБД



Показывает статус и взаимосвязи сервисов СУБД
Акцентирует внимание на неисправных объектах
Автоматизирует операцию switchover сервиса СУБD



спектр

ОБЪЕКТЫ УПРАВЛЕНИЯ

Кластеры

Узлы

Сервисы СУБД

ИСТОРИЯ

Операции

ДОСТУПЕН HA-alt82-2302

ОбзорКонфигурацияКонфигурация СУБДТопологияЭкземпляры сервисаМониторинг

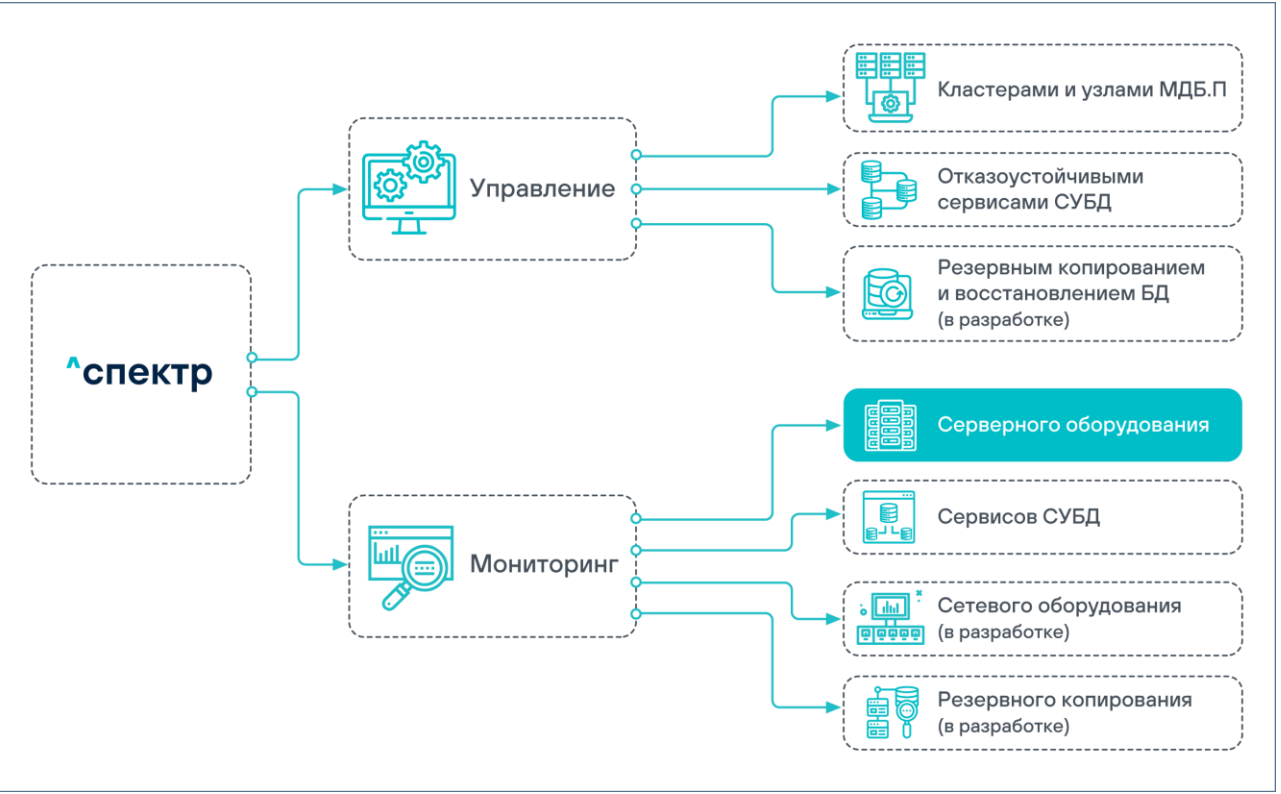
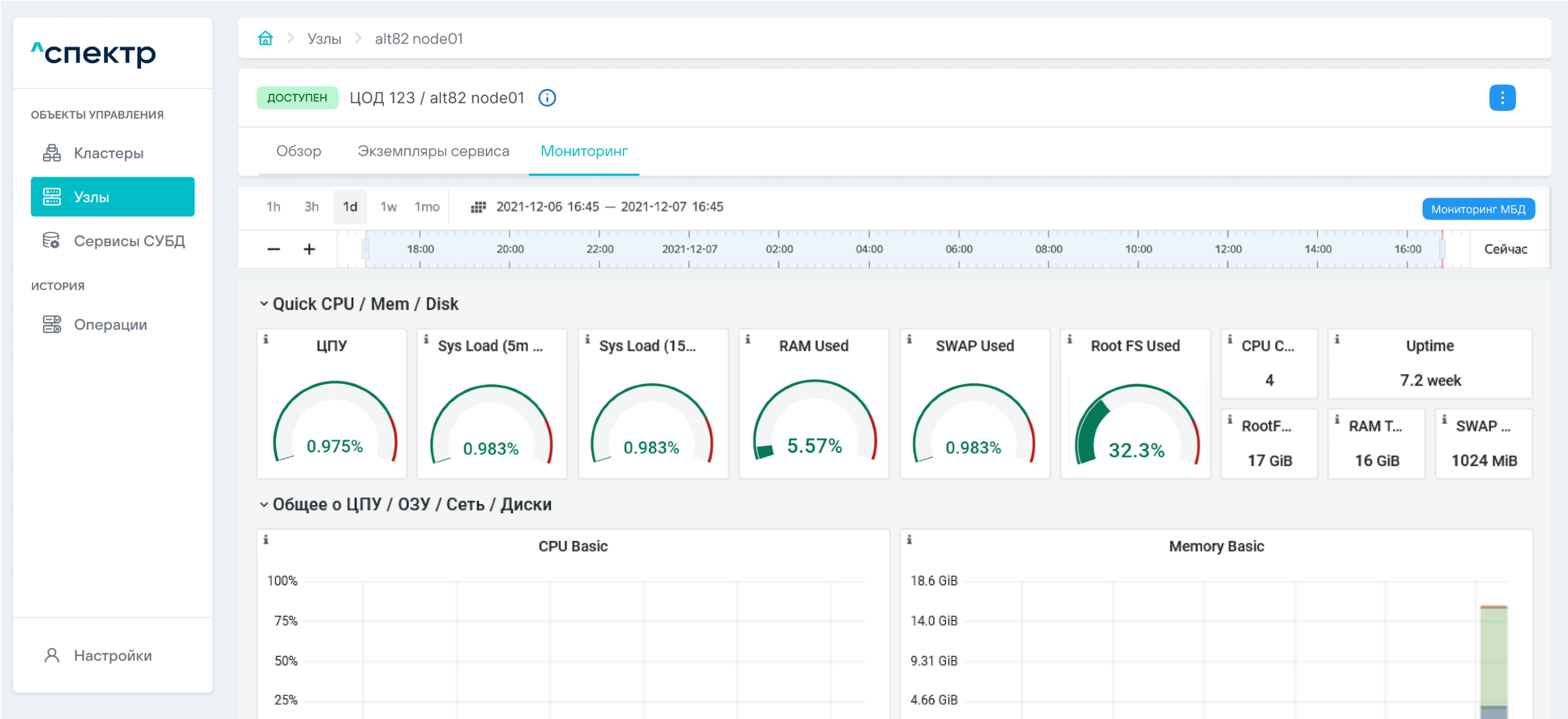
О сервисе

ПлощадкаЦОД 123ДОСТУПЕНВерсия PostgresPostgres Pro Standard 14.2.1

Имя кластера:alt82ДОСТУПЕНUptime мастера254d

Role	IP	Узел	Статус	Flush lag	Replay lag	Active conn.	Max conn.	TPS	QTS
PRI	192.168.20.5:5432	alt82 node01	ДОСТУПЕН	--	--	35	3000	85	92
HS:async	192.168.20.5:8080	alt82 node03	ДОСТУПЕН	1с	10с	1	3000	24	34
HS:sync	255.25.0.1:5432	alt82 node02	ДОСТУПЕН	0с	0с	46	3000	79	85

Мониторинг серверного оборудования



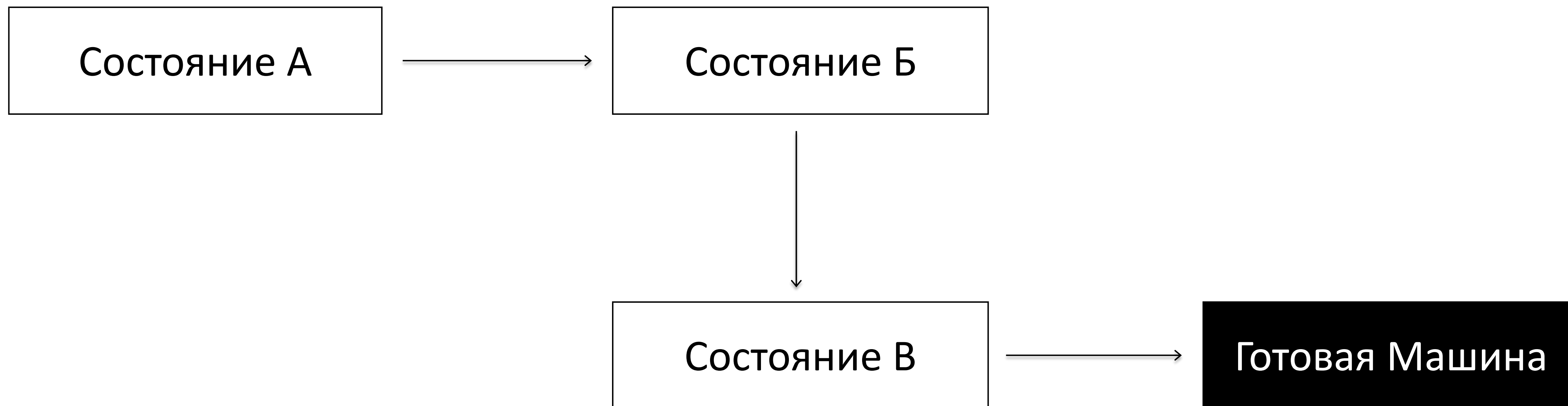
Мониторинг СУБД



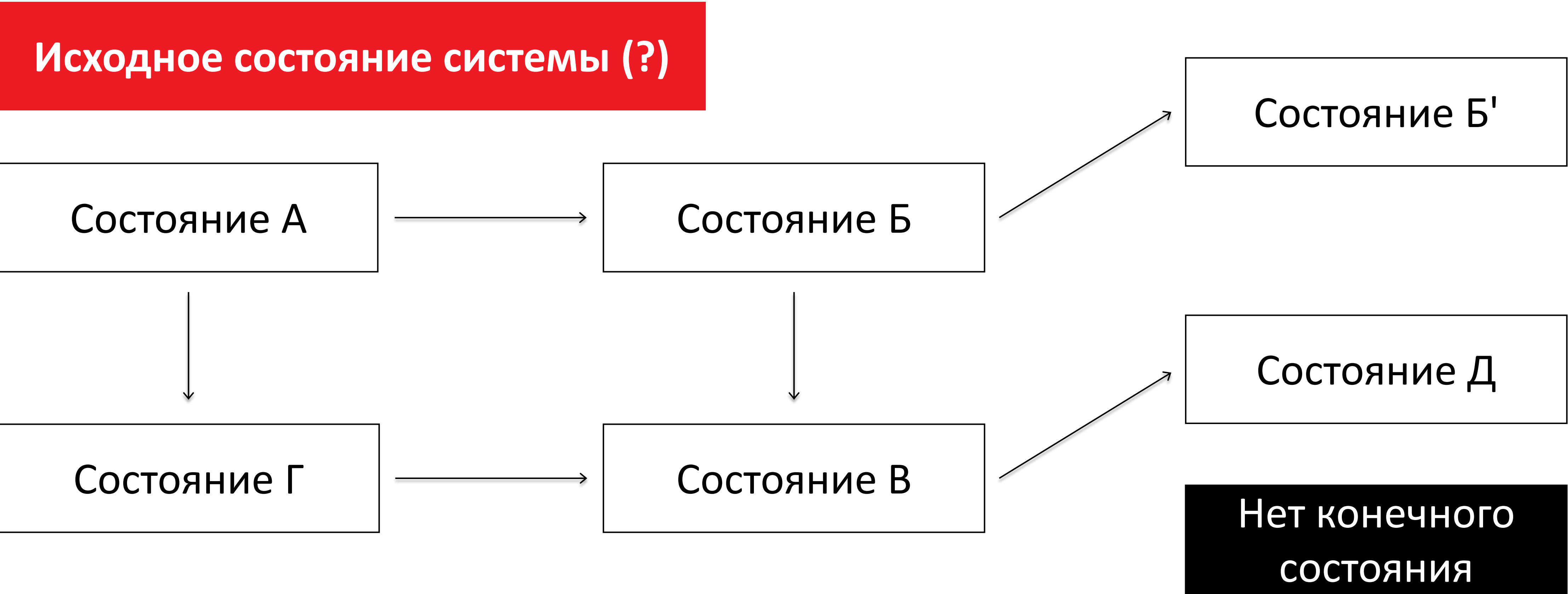
Спектр. О пользе и вреде копирования



Исходное состояние системы



Спектр. О пользе и вреде копирования

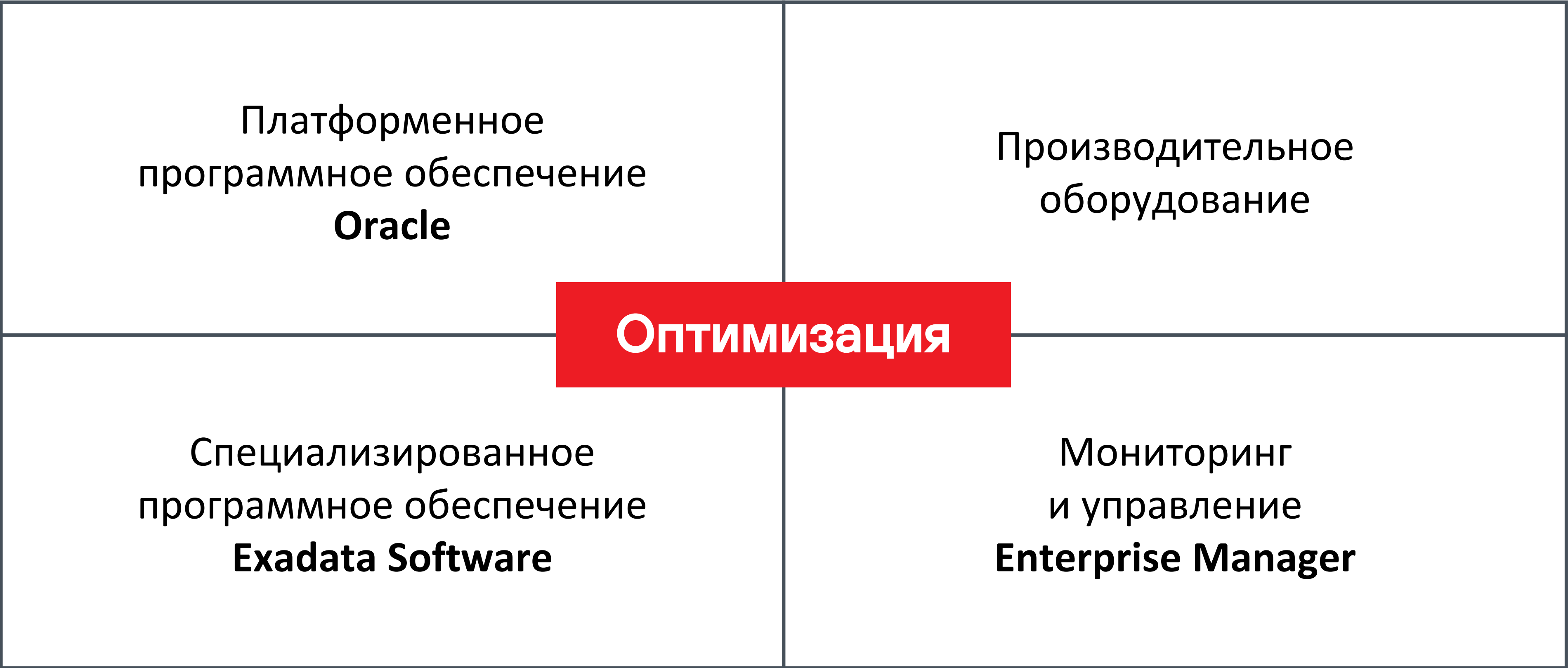


Спектр. Непрошенные советы

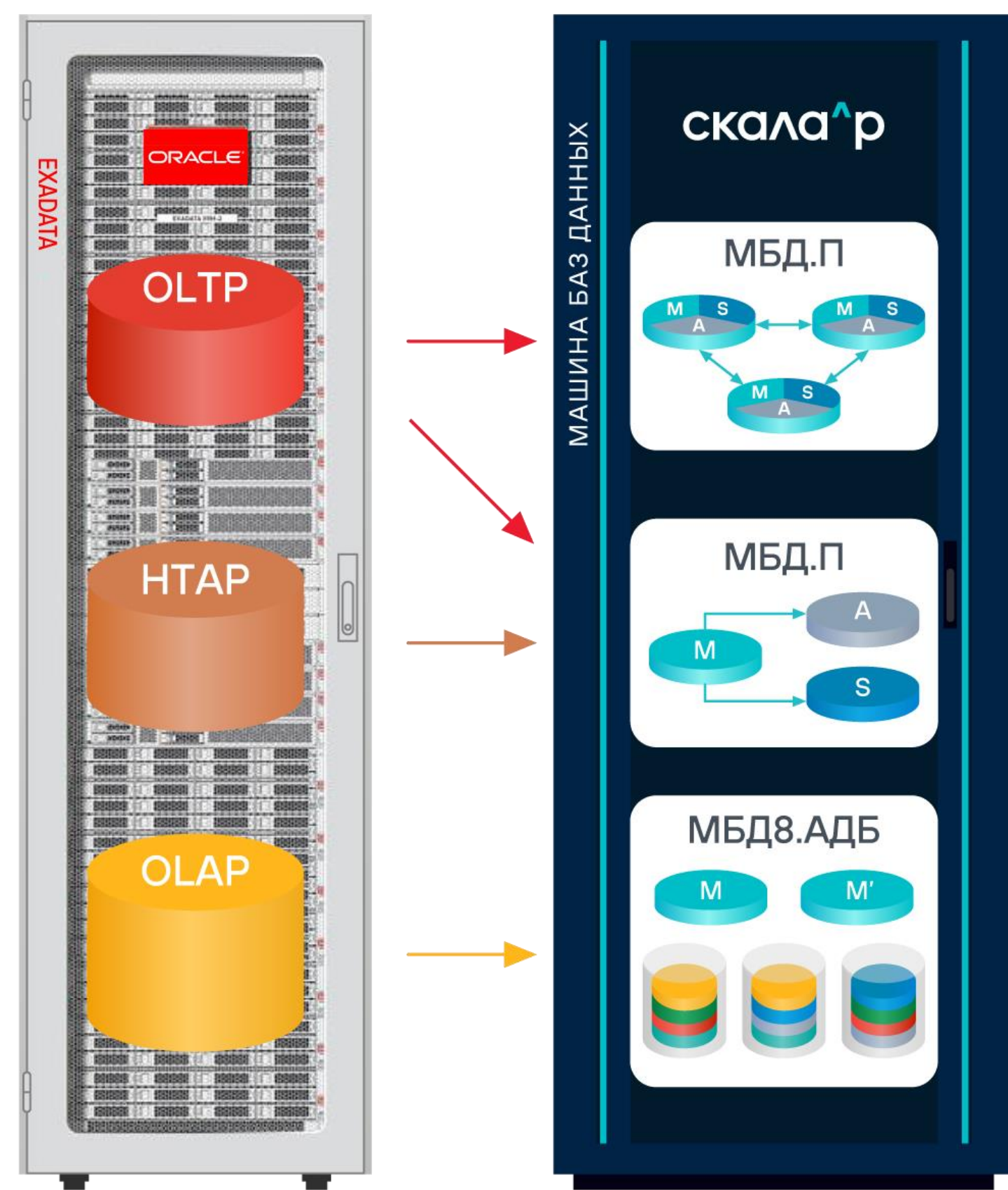


```
<crm_mon version="2.1.1-alt1">
  <summary>
  <nodes>
  <resources>
    <clone id="HA-1701-pg11-a9-1-clone" multi state="true" unique="false"
      failure ignored="false">
    <group id="HA-1701-pg11-a9-1-master-group" number_resources="2" manage
    <resource id="HA-1701-pg11-a9-1-master_ip" resource_agent="ocf:heart
      "false" blocked="false" managed="true" failed="false" failure_ignore
    <node name="1701-pg11-a9-1-01" id="1" cached="true"/>
    </resource>
    <resource id="HA-1701-pg11-a9-1-master repl ip" resource_agent="ocf
      orphaned="false" blocked="false" managed="true" failed="false" failu
    <node name="1701-pg11-a9-1-01" id="1" cached="true"/>
    </resource>
    </group>
    <resource id="HA-1701-pg11-a9-1-sync_ip" resource agent="ocf:heartbea
      "false" blocked="false" managed="true" failed="false" failure ignored=
    <resource id="HA-1701-pg11-a9-1-async_ip" resource agent="ocf:heartbea
      "false" blocked="false" managed="true" failed="false" failure ignored=
```

Best Practice в создании ПАК

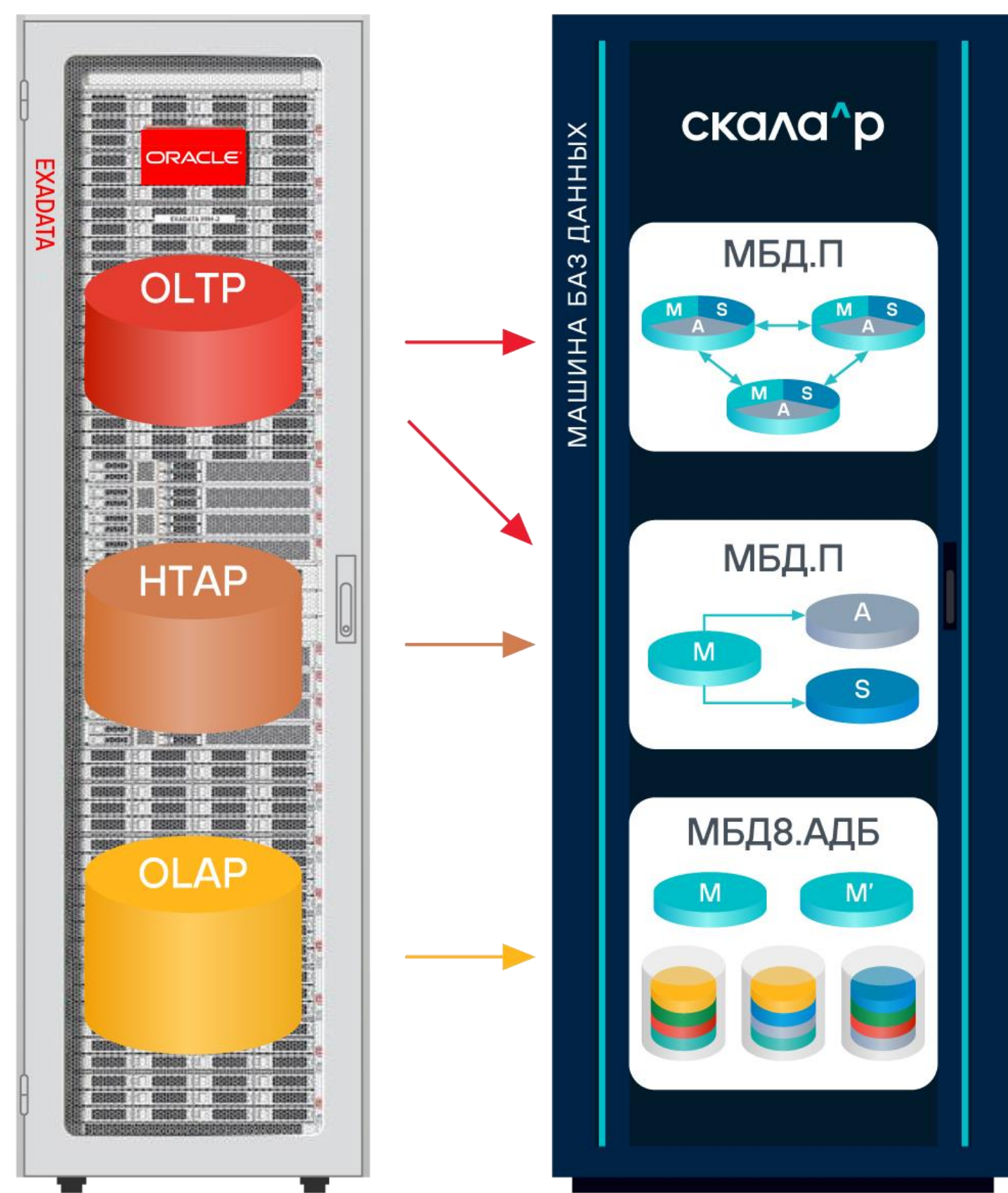


Миграция Oracle Exadata -> Скала^р МБД.П



- Скала^р МБД.П в инфраструктуре заказчика
- Миграция данных (Ora2PG)
- Profit!!!

Миграция Oracle Exadata -> Скала^р МБД.П



- Скала^р МБД.П в инфраструктуре заказчика — ✓
- Миграция данных (Ora2PG) — ✓
- ~~Profit!!!~~

Замена СУБД в ИС



Можете привести примеры, где несложно импортозаместиться?

...

PostgreSQL, мне кажется, тоже хорошая история. Но не везде и не всегда.
Oracle так просто не заменишь, но во многих случаях можно.

...А в ФНС возможно заменить?

В АИС «Налог-3» — невозможно.

Что для этого с Postgres нужно сделать?

Ничего. Невозможно.

Тотальная зависимость от Oracle?

Да, надо просто написать новую систему на другой архитектуре.

4-ю версию АИС «Налог»?

Да.



Замена СУБД в ИС




Как мы переписывали бизнес-логику высоконагруженного приложения на PLPG/SQL

Базы данных и системы хранения

Миграция / Хранимые процедуры

#Миграция / Хранимые процедуры

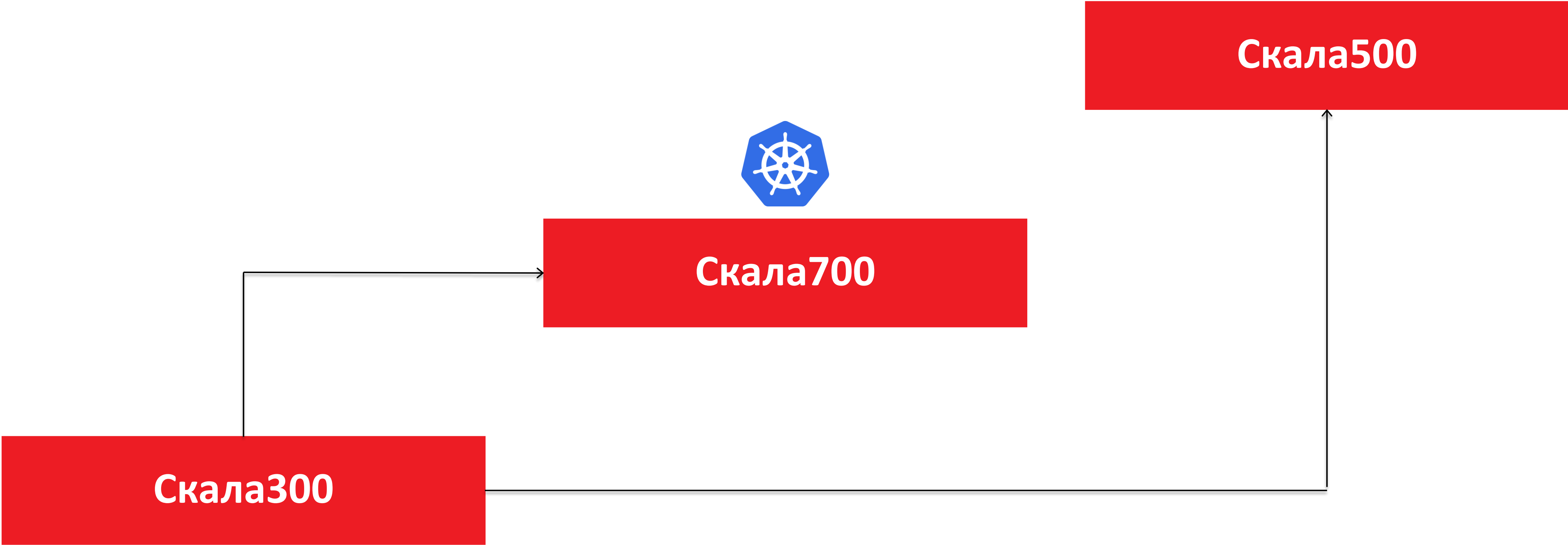
25 ноября, 10:00, Зал «h3: Яндекс трек» 

 Доклад принят в программу конференции

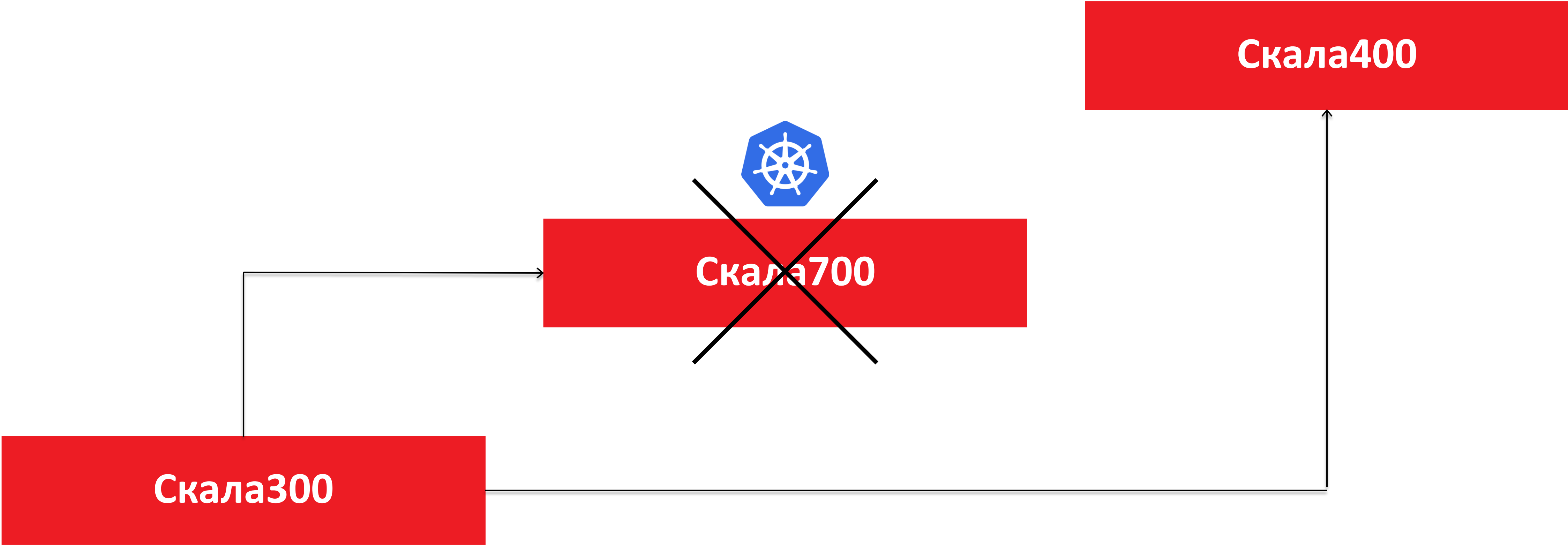
Мнение Программного комитета о докладе

Миллионы наших сограждан оформляют больничные, а ИТ-система в ФСС, которая обслуживает этот процесс, переехала из Oracle в Postgres. В докладе будет все об особенностях этого увлекательного процесса.

Скала300 -> Скала500 -> Скала700



Скала300 -> Скала500 -> Скала700

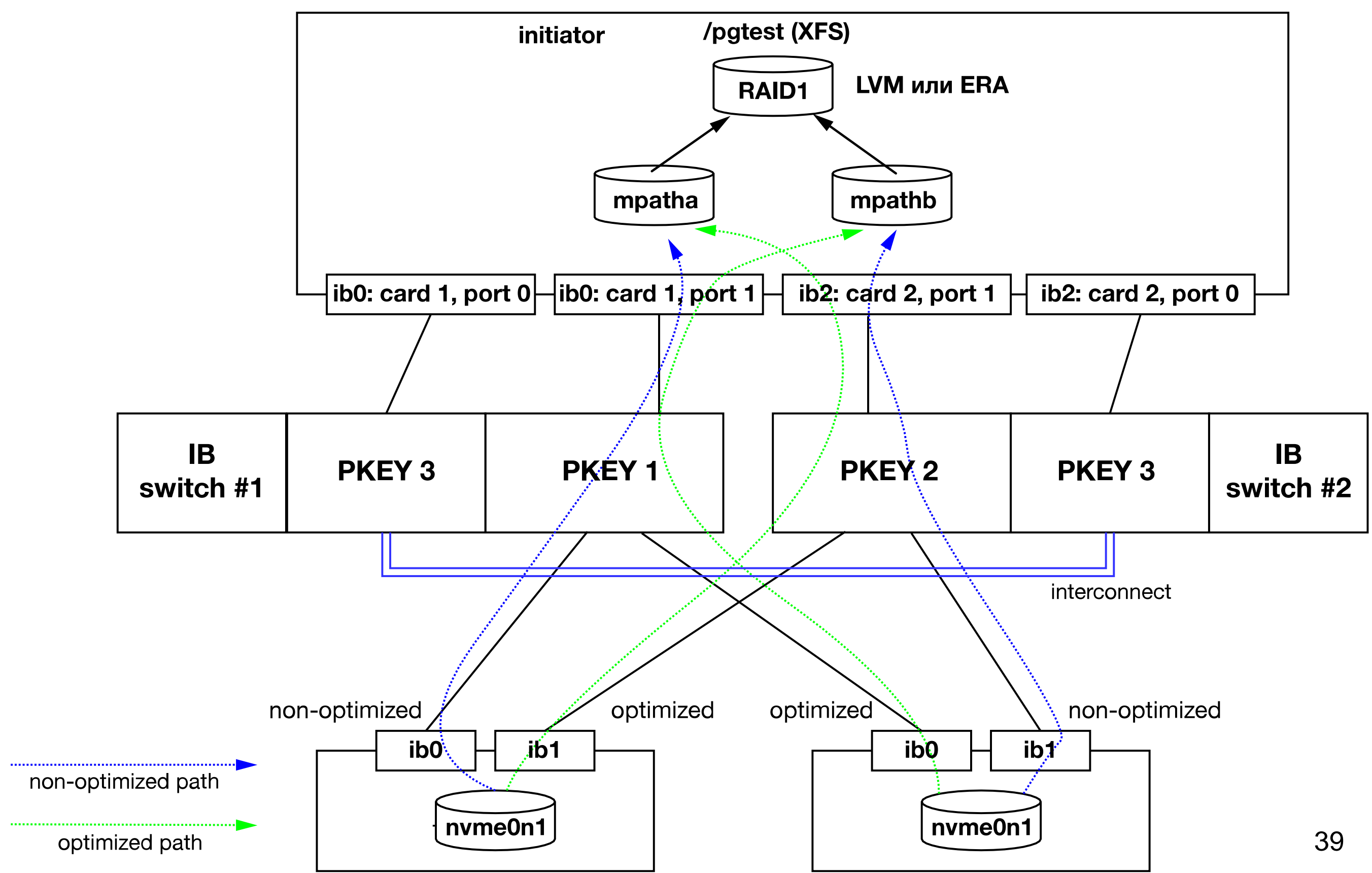


Скала500 NVMe over Fabric



Скала300

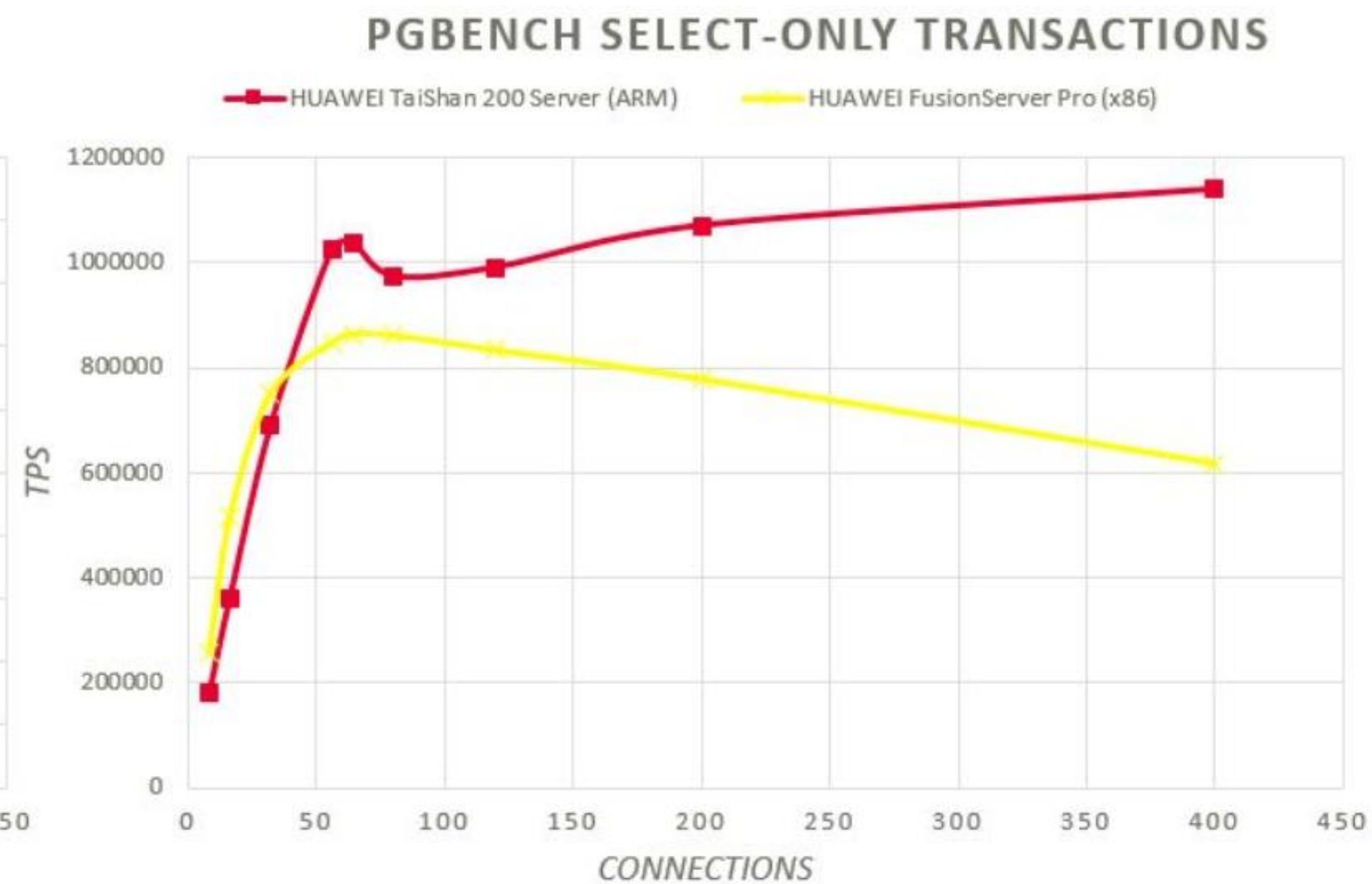
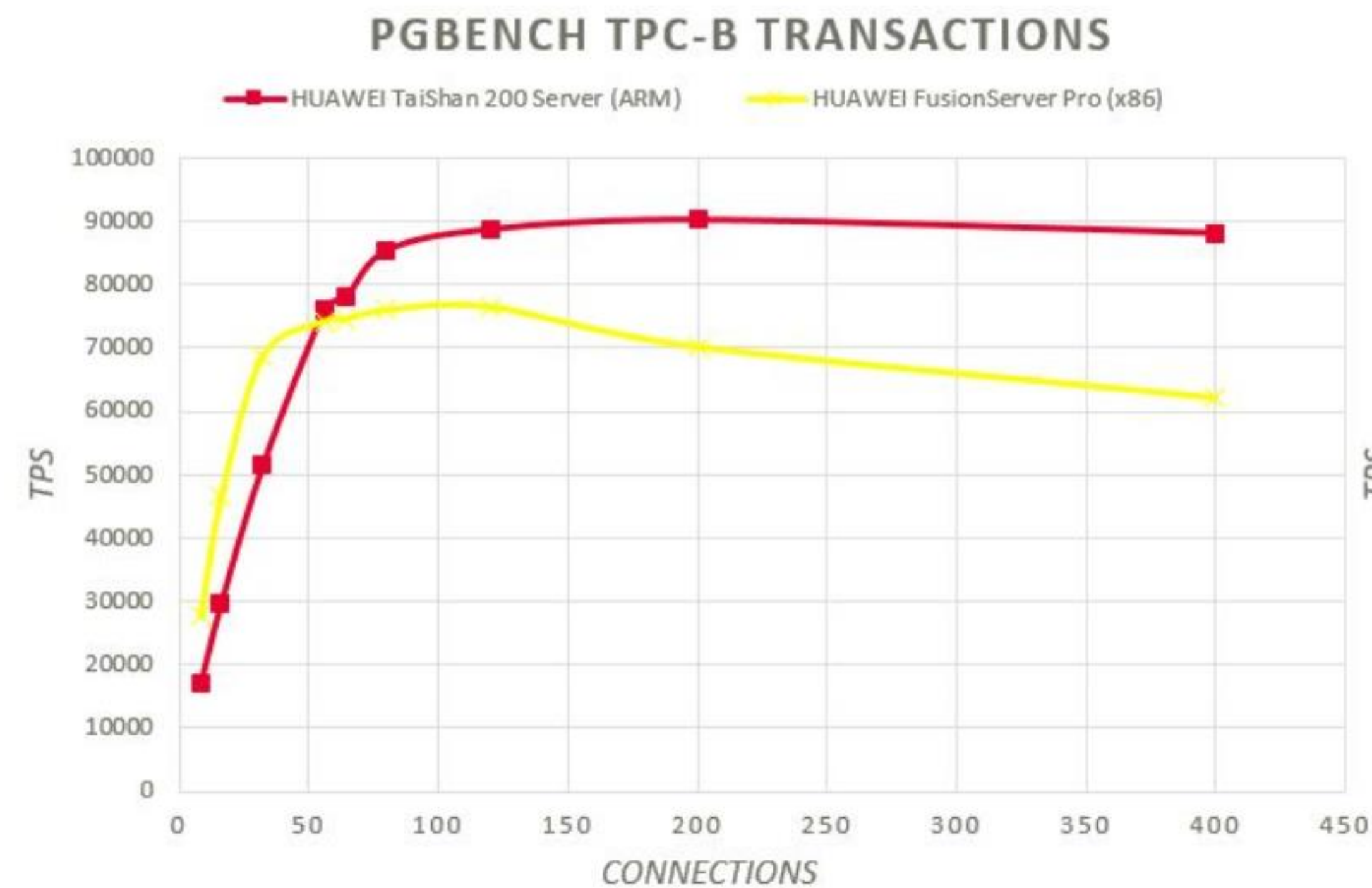
Скала500



Про CPU: Xeon & Epic -> Эльбрус и Байкал -> ARM64

Архитектура	Модель	Доступность	Примечание
x86-64	Intel Xeon/AMD Epic	+/-	Полностью устраивают
VLIW	Эльбрус	☹	Производительность?
ARM64	Байкал-M	☹	Не успели попробовать
ARM64	«Китайские» ARM	+/-	Производительность?
RISC-V	Байкал-S	?	Ждём
RISC-V	Решения от Yadro	?	Ждём

Тесты на Huawei TaiShen



Тесты на Huawei TaiShen



Пример конфигурации сервера БД PostgreSQL (postgresql.conf):

```
max_connections = 1024
shared_buffers = 390GB
max_prepared_transactions = 2048
huge_pages = try
work_mem = 1GB
maintenance_work_mem = 2GB
dynamic_shared_memory_type = posix
```

Это не хорошо и не плохо, просто результаты нерелевантны для наших условий!!!

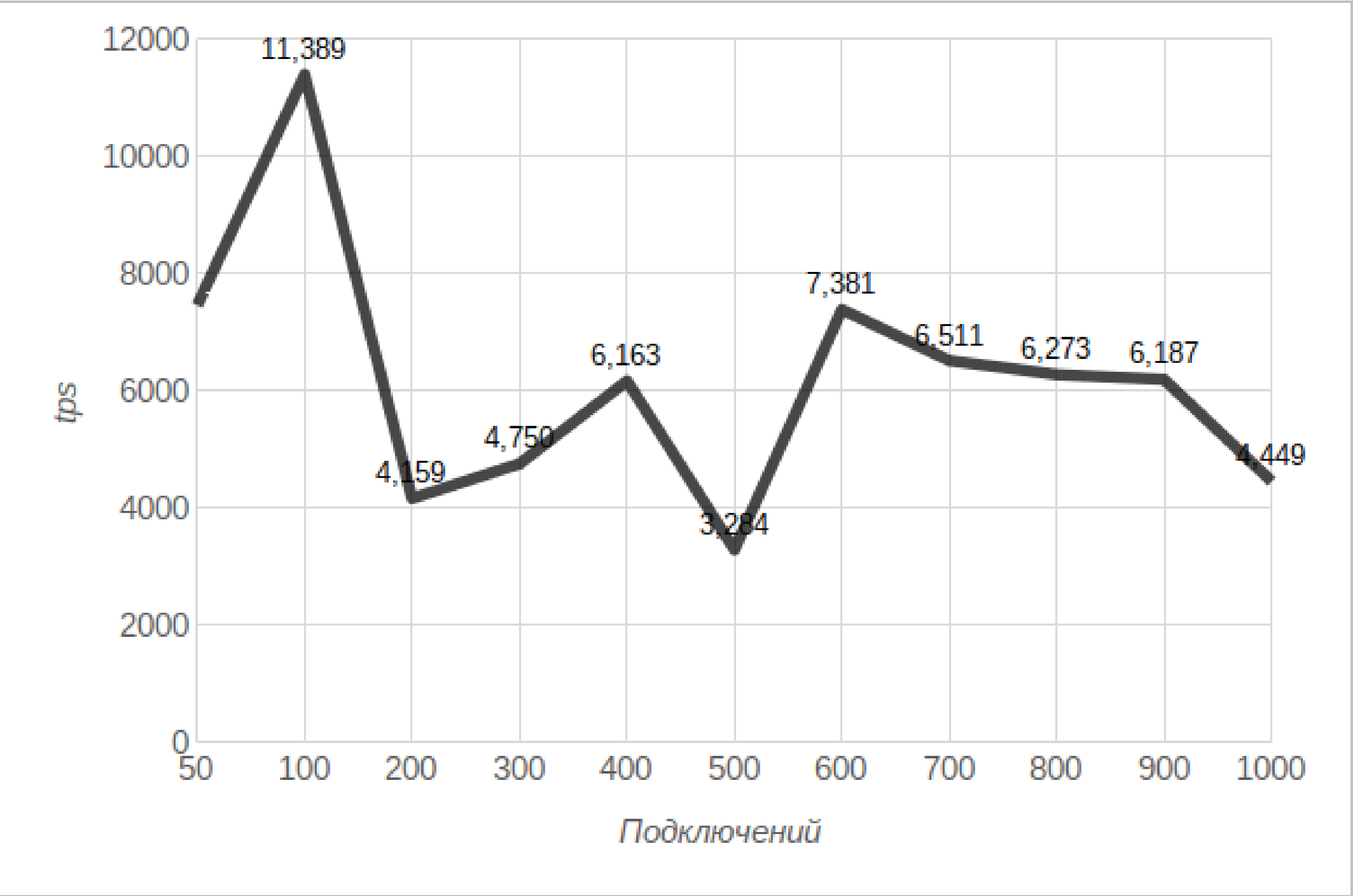
Методика:

- 1) Создаем тестовую БД следующей командой (размер около 370 ГБ).

```
pgbench -i -s 25000
```

- 2) Выполняем несколько SQL-запросов для разогрева кэша БД:

Тестирование TaiShen в IBS Интерлаб



Знаешь Скала^р?
Отметься!





Ответвление отPostgreSQL 9.2		
C++ вместо C	Thread-per-connection вместо process-per-connection	
Основная особенность — трёхдвижковость		
Строчный	Резидентный	Столбцовый
Специфика по работе с NUMA (Куньлун, Куньпен)		
Модифицированное журналирование		
Многопоточное применение журнала предзаписи	Журналы отмены (практически REDO LOG у Oracle!!!)	



СКАЛА^р
отсекая лишнее

Спасибо!

Константин Аристов

Скала^р, техлид

Karistov@skala-r.ru

@A_K_M_74



HighLoad ++
2022

Яндекс

Оценка

